# Feedback and therapist effects in the context of treatment outcome and treatment length

Wolfgang Lutz, Julian Rubel, Anne-Katharina Schiefele, Dirk Zimmermann, Jan Rasmus Böhnke & Werner W. Wittmann

Published online: 28 Jul 2015.

Submit your article to this journal 🗗

Article views: 89

View related articles 🗗

View Crossmark data 🗗

Citing articles: 1 View citing articles 🗗

**EMPIRICAL PAPER**

# Feedback and therapist effects in the context of treatment outcome and treatment length

WOLFGANG LUTZ[1], JULIAN RUBEL[1], ANNE-KATHARINA SCHIEFELE[1],
DIRK ZIMMERMANN[1], JAN RASMUS BÖHNKE[2], & WERNER W. WITTMANN[3]

[1]*Department of Clinical Psychology and Psychotherapy, University of Trier, Trier, Germany;* [2]*Hull York Medical School & Department of Health Sciences, University of York, York, UK* & [3]*Otto-Selz-Institute for Applied Psychology, University of Mannheim, Mannheim, Germany*

**Abstract**
**Objective:** This study estimates feedback and therapist effects and tests the predictive value of therapists' and patient attitudes toward psychometric feedback for treatment outcome and length. **Methods:** Data of 349 outpatients and 44 therapists in private practices were used. Separate multilevel analyses were conducted to estimate predictors and feedback and therapist effects. **Results:** Around 5.88% of the variability in treatment outcome and 8.89% in treatment length were attributed to therapists. There was no relationship between the average effectiveness of therapists and the average length of their treatments. Initial impairment, early alliance, number of diagnoses, feedback as well as therapists' and patients' attitudes toward feedback were significant predictors of treatment outcome. Treatments tended to be longer for patients with a higher number of approved sessions by the insurance company, with higher levels of interpersonal distress at intake, and for those who developed negatively (negative feedback) over the course of their treatment. **Conclusions:** Therapist effects on treatment outcome and treatment length in routine care seem to be relevant predictors in the context of feedback studies. Therapists' attitudes toward and use of feedback as well as patients' attitudes toward feedback should be further investigated in future research on psychometric feedback.

**Keywords:** attitudes towards feedback; feedback effects; patient-focussed (psychotherapy) research; patient reported outcome; therapist effects

## Introduction

Patient-focused research (PFR) represents one important research strategy to address the widely complained scientist–practitioner gap in providing feedback tools to implement research results "in real time" into clinical practice (e.g., Castonguay, Barkham, Lutz, & McAleavy, 2013; Howard, Moras, Brill, Martinovich, & Lutz, 1996; Lambert, 2007). Central to PFR is the continuous monitoring of patient outcomes over the course of treatment and a psychometric feedback of this information to therapists (and potentially patients). These assessment results are used to evaluate the current state and to predict the future course of an individual patient and can signal early, potentially negative developments and initiate a change in treatment strategy (e.g., Hannan et al., 2005; Meehl, 1954). Research suggests that this enhanced emphasis on outcome orientation and psychometric feedback is a promising path for a further improvement of the effectiveness of psychotherapy, especially for those patients showing negative developments early in treatment (Carlier et al., 2012; Shimokawa, Lambert, & Smart, 2010). In addition, some studies showed continuous feedback could support an optimized resource allocation since patients who show negative developments stay longer in treatment than patients who show positive developments, whereas patients with positive developments (on track) stay shorter in

treatment (e.g., Lambert & Shimokawa, 2011; Lambert et al., 2003). A continuous outcome monitoring and feedback could be also a promising path to differentiate patients who need more sessions of treatment from those who already profit from smaller amounts (e.g., Lutz, Ehrlich, et al., 2013; Stulz, Lutz, Leach, Lucock, & Barkham, 2007).

However, early feedback studies were conducted within settings in which relatively short treatments were provided to only moderately impaired patients (e.g., college counseling centers; Newnham & Page, 2010; Lutz et al., 2014; Poston & Hanson, 2010; Shimokawa et al., 2010). Recent studies investigated feedback effects in more disturbed outpatients (De Jong, van Sluis, Nugter, Heiser, & Spinhoven, 2012; Simon, Lambert, Harris, Busath, & Vazquez, 2012), in psychosomatic in-patients (Byrne, Hooke, Newnham, & Page, 2012; Probst et al., 2013), in patients with eating disorders (Simon et al., 2013), and for patients receiving long-term treatments (≥35 weeks; De Jong et al., 2014). These more recent investigations generally showed less pronounced feedback effects. Different explanations have been discussed for this reduced effect (Riemer & Bickman, 2011; Simon et al., 2012). A closer look revealed that feedback is not uniformly effective for every patient and therapist. De Jong et al. (2012), for example, found substantial differences between therapists regarding their use of feedback. Having a higher commitment to use the feedback as well as being female showed to be significantly associated with a higher probability to make use of the feedback for the ongoing treatment. Therapists in turn who reported to use the feedback showed to be more effective for negatively developing patients. Additionally, patients of therapists who were more committed to use the feedback at the beginning of the study showed a faster treatment response. Likewise, Simon et al. (2012) reported that only 50% of the therapists in their study were able to use the feedback to substantially improve client' outcomes. For the other half of therapists it made no difference whether or not they got feedback about the progress of their patients. Until now, only these two studies investigated therapist differences and their effects on outcome in the context of feedback studies. But, these investigations have been limited to treatment outcome and did not examine the influence of these variables on treatment length.

Related to the differences between therapists reported in these feedback studies is the discussion about therapist differences in their ability to successfully provide psychological interventions (Baldwin & Imel, 2013; Saxon & Barkham, 2012). Theoretically, it is quite intuitive that therapists differ with regard to their average effectiveness, their ability to form a sustainable relationship with their clients, their reasoning regarding how much therapy is enough, and other process-relevant variables. Research suggests that about 5–8% of the variability in outcomes is due to therapist differences and depends in part on differences in the study designs and the severity of the patients' impairment (Baldwin & Imel, 2013; Crits-Christoph et al., 1991; Saxon & Barkham, 2012). It seems that the effect increases for more severely impaired patients and is higher in naturalistic compared to controlled samples (Lutz, Leon, Martinovich, Lyons, & Stiles, 2007; Saxon & Barkham, 2012).

Most research from naturalistic samples has focused on therapist effects on outcome and findings on treatment length are sparse. Randomized controlled trials (RCTs) are in this regard only of limited usefulness, since they usually have a treatment protocol with a fixed number of treatment sessions. Therapist effects on treatment length could unveil systematic differences in treatment lengths which are due to the therapists' individual concepts of how much treatment is enough or therapists' ability to maintain a treatment and not letting patients dropout. The identification of therapists who generally provide longer treatments is also an important issue in the context of scarce financial resources in mental health-care settings. These therapists might profit from tools helping to differentiate patients who need more sessions from patients who need fewer sessions. Only one feedback study systematically investigated therapist effects on treatment length. Okiishi, Lambert, Nielsen, and Ogles (2003) found that patients seen by the best three therapists in their study not only experienced on average more change in OQ-45 scores from pre- to post-treatment but also stayed shorter in treatment than those patients seen by the three worst therapists. However, systematic differences between therapists' average treatment lengths have not been investigated independently from treatment outcome.

Therefore, in the current investigation first the question about how psychometric feedback, therapists' as well as patients' attitudes toward feedback and therapist differences influence treatment outcome and treatment length is addressed when controlling for several additional patient intake characteristics. The second question addresses whether there is an association between therapists' average treatment length and therapists' average treatment outcome.

## Methods

### The "Techniker Krankenkasse (TK) Project"

In collaboration with the TK, a German health insurance company, a quality monitoring study was

conducted, which included decision rules and feedback tools. The *TK project* was the first study to evaluate the practical feasibility of a quality assurance and feedback system in private practices within the German health insurance system (Lutz, Wittmann, Böhnke, Rubel, & Steffanowski, 2013; Strauss et al., 2015; Wittmann et al., 2011). The main goal was to test whether a quality management strategy using feedback was feasible and associated with better outcomes compared to therapies, which were subject to the traditional peer review system of quality assurance (see also Lutz, Wittmann et al., 2013; Strauss et al., 2015; Wittmann et al., 2011). This study was designed as a cluster-RCT involving psychotherapists from three different therapeutic approaches: cognitive behavioral therapy (CBT), psychodynamic therapy and psychoanalysis.[1] The participating therapists were randomly allocated to (i) the traditional case report model of quality assurance including psychometric assessments at intake, termination and follow-up (control, CG) and (ii) an alternative approach to quality management including continuous psychometric testing and feedback (intervention, IG). Depending on how long patients were treated, self-reports were obtained after session 10, 20, 40, 55, and 75 in CBT and after session 10, 20, 45, 55, 75, and 95 for psychodynamic treatments. Each assessment was fed back to the therapists in the IG within a few days whereas no information was provided to the therapists of the CG. To increase the usefulness of the feedback, therapists from the IG were provided with rationally derived decision rules about their patients' progress based on an extension of clinically significant change criteria (for more details on the design of the feedback see Lutz, Böhnke, & Köck, 2011). Therapists who participated in the IG were provided with a simplified approval procedure. Instead of the usual 25 sessions more sessions were approved for treatments in the intervention group (CBT: 45; PD: 50). This modification to the usual application procedure resulted in a higher number of approved sessions in the IG compared to the CG (Table I).

### Study Sample

In the current study, a subsample of the TK health insurance sample ($N = 751$, 177 therapists) was used, which included only therapists who treated at least 5 patients. This selection was necessary to reach a minimal precision for the estimates of within- and between-therapist variability for determining the therapist effects (see also Baldwin et al., 2011). Furthermore, only patients were included who provided a pre- and post-score in the Brief Symptom Inventory (BSI; Franke, 2000) as well as additional

information on diagnoses, early alliance, and feedback (IG). Thus, the analysis sample included 44 therapists and 349 patients in which each therapist treated between 5 and 18 patients ($M = 9.25$, $SD = 3.84$, median = 8). A standardized assessment procedure including a structured diagnostic interview (related to ICD-10 criteria; Hiller, Zaudig, & Mombour, 2004) was conducted with the patients from the IG. In the CG diagnoses were based on clinical judgments. The primary axis one diagnoses in the study sample were distributed as follows: 39.0% of the patients had a major depressive disorder ($n = 136$), 9.2% ($n = 32$) had a dysthymic disorder, 20.1% ($n = 70$) had an adjustment disorder, 18.9% ($n = 66$) had an anxiety disorder, 2.0% ($n = 7$) had an eating disorder, 8.6% ($n = 30$) had other diagnoses, and 2.3% ($n = 8$) were not diagnosed with an axis one disorder. Additionally, about 10.3% ($n = 36$) of the patients were diagnosed with a personality disorder. The number of diagnoses in the analyzed sample ranges from one to four with an average value of $M = 1.52$. Further demographic characteristics for patients and therapists are provided in Tables I and II, respectively.

### Instruments and Data Collection

**Brief Symptom Inventory (BSI).** Symptom severity was measured using the BSI (Franke, 2000; German translation of Derogatis, 1975), a 53-item self-report inventory, which asks about physical and psychological symptoms within the last week. It is the short form of the Symptom Check-List-90-Revised (SCL-90-R; Derogatis, 1975). Item responses are based on a 5-point Likert scale ranging from 0 ("not at all") to 4 ("extremely"). As the primary therapy outcome indicator, the Global Severity Index (GSI) was computed for pre- and post-treatment by averaging all BSI items. For this index an internal consistency of $\alpha = .92$ and a retest-reliability of $r_{tt} = .90$ are reported (Franke, 2000).

**Inventory of Interpersonal Problems (IIP-D).** Initial interpersonal distress was assessed with the IIP-D (Horowitz, Strauß, & Kordy, 2000). For the present project the 64-item version of the IIP-D was used, which measures interpersonal problems regarding behavior, thoughts, and emotions. Item responses are based on a 5-point Likert scale ranging from 0 ("not") to 4 ("very"). For the overall mean score an internal consistency of $\alpha = .94$ and a retest-reliability (10 weeks) of $r_{tt} = .98$ are reported (Horowitz et al., 2000).

**Penn Helping Alliance Questionnaire (HAQ).** Therapeutic alliance was assessed from a patient's

Table I. Descriptive patient characteristics (Mean [*SD*]) and 95% CI for the completer sample, the study sample, and the sample of excluded patients grouped by the treatment condition and on average.

| | Completer sample | | Study sample | | Excluded sample | |
|---|---|---|---|---|---|---|
| | IG (*n* = 507) | CG (*n* = 244) | IG (*n* = 268) | CG (*n* = 81) | IG (*n* = 239) | CG (*n* = 163) |
| Age | 44.4 (11.4) | 45.6 (11.1) | 44.5 (11.3) | 46.0 (10.5) | 44.3 (11.4) | 45.3 (11.5) |
| | 44.8 [43.97, 45.57] | | 44.8 [43.65, 46.00] | | 44.8 [43.59, 45.83] | |
| Female | 67.3% | 68.0% | 64.2% | 65.4% | 70.7% | 69.3% |
| | 67.5% [64.1%, 70.8%] | | 64.5% [59.3%, 69.3%] | | 70.2% [65.5%, 74.4%] | |
| Number of diagnoses | 1.61 (.75) | 1.38 (.65) | 1.56 (.73) | 1.36 (.68) | 1.69 (.78) | 1.40 (.64) |
| | 1.53 [1.48, 1.59] | | 1.52 [1.44, 1.59] | | 1.55 [1.47, 1.63] | |
| Approved sessions[a] | 50.96 (16.4) | 35.63 (18.4) | 48.35 (9.4) | 37.32 (22.7) | 54.1 (21.5) | 34.8 (15.7) |
| | 46.08 [44.71, 47.44] | | 45.85 [44.32, 47.38] | | 46.29 [44.08, 48.50] | |
| Number of sessions | 42.71 (21.4) | 36.18 (19.0) | 40.94 (19.0) | 35.99 (19.8) | 44.8 (23.8) | 36.3 (18.6) |
| | 40.63 [39.09, 42.18] | | 39.82 [37.76, 41.88] | | 41.38 [39.10, 43.66] | |
| $GSI_{pre}$ | 1.25 (.67) | 1.14 (.63) | 1.26 (.67) | 1.09 (.62) | 1.24 (.68) | 1.16 (.63) |
| | 1.21 [1.16, 1.26] | | 1.22 [1.15, 1.29] | | 1.20 [1.14, 1.27] | |
| $GSI_{post}$ | 0.62 (.56) | 0.55 (.47) | .63 (.58) | .50 (.38) | .60 (.55) | .57 (.51) |
| | .59 [.56, .63] | | .60 [.54, .66] | | .59 [.54, .64] | |
| $d_{GSI}$ | 0.96 (.93) | 0.90 (.97) | 0.96 (.91) | 0.90 (1.00) | .96 (.95) | .89 (.97) |
| | .94 [.87, 1.01] | | .95 [.85, 1.04] | | .93 [.84, 1.03] | |
| $d_{GSIadjusted}$ | 0.89 (.77) | 0.98 (.77) | 0.92 (.72) | 1.02 (.72) | .96 (.73) | .96 (.73) |
| | .94 [.88, .99] | | .97 [.88, 1.06] | | .96 [.89, 1.03] | |
| $HAQ_{pre}$ | 53.8 (7.1) | 56.2 (6.9) | 53.4 (7.5) | 55.8 (6.4) | 54.3 (6.4) | 56.4 (7.1) |
| | 54.58 [54.06, 55.10] | | 53.97 [53.21, 54.74] | | 55.16 [54.46, 55.86] | |
| $IIP_{pre}$ | 1.55 (.55) | 1.49 (.55) | 1.53 (.54) | 1.46 (.55) | 1.58 (.56) | 1.51 (.55) |
| | 1.53 [1.49, 1.57] | | 1.53 [1.46, 1.57] | | 1.55 [1.50, 1.61] | |

[a]Differences between control and intervention groups in terms of approved and number of sessions are due to the design of the original study (see also Strauss et al., 2015, and Methods section).

perspective with the HAQ after the first and after every 10th session (Alexander & Luborsky, 1986; German translation of Bassler, Potratz, & Krauthauser, 1995). The HAQ is an 11-item self-report questionnaire with a 6-point Likert scale ranging from 1 ("strongly agree") to 6 ("strongly disagree") so that lower mean scores imply a better alliance.

**Feedback.** Feedback in the TK project was based on three self-report questionnaires: the BSI (Franke, 2000) to cover general psychological distress, the IIP (Horowitz et al., 2000) to assess interpersonal impairment, and a disorder-specific instrument taking into account the main diagnosis of each patient (Beck Depression Inventory, Beck, Steer, & Carbin 1988; AKV – anxiety disorders, Ehlers & Margraf, 2001; Eating Disorder Inventory, Garner, 1991; HZI-K – compulsive disorders, Klepsch, Zaworka, Hand, Lünenschloss, & Jauering, 1993; SOMS – somatoform disorders, Rief, Hiller, & Heuser, 1997). Psychometric feedback after each assessment was generated based on clinical significant change calculations in relation to intake scores (Jacobson & Truax, 1991). The results of each instrument were integrated into an overall evaluation of patient progress, which described progress by one of three feedback reports: (i) "overall negative change", (ii) "no reliable change so far", and (iii) "good progress" indicating that the

patient changed positively compared to the initial assessment. To be able to include patients from the control group in the analyses, feedback was dummy coded ("one negative feedback over the course of the treatment" and "no negative feedback over the course of the treatment") with the control group as reference category. For all patients from the IG group feedback after the 10th session was available. Due to varying treatment lengths, after the second assessment (session 20) feedback data were available for 93.3% of the patients in the IG, for 78.7% after the third, for 52.6% after the fourth, for 26.1% after the fifth, and for 9% after the sixth assessment.

**Assessment of patients' attitudes toward routine outcome monitoring.** Patients' attitudes toward routine outcome monitoring were assessed at the end of their therapy. An index of patients' attitudes toward feedback was composed of six items. In five of these items patients report their perception of the continuous outcome assessment procedure (e.g., "I find it important to monitor the results of psychotherapeutic treatments.") on a 5-point Likert scale ranging from 1 ("agree completely") to 5 ("disagree completely"). Additionally the index comprises one global item measuring patients' overall satisfactions with the project ("All together, how satisfied are you with the project?") on a 5-point

Likert scale ranging from 1 ("completely satisfied") to 5 ("dissatisfied"). The average of these six items was used as an indicator of patients' attitudes toward feedback. Thus, lower scores on that index imply more positive attitudes toward feedback. To be able to include patients from the CG and patients with missing values (to maintain statistical power for our analysis), this index was effect-coded: "very positive attitude" ($n = 44$; 12.6%; upper third), "average and negative attitude" ($n = 90$; 25.8%; lower two thirds), "missing" ($n = 134$; 38.4%), and "control group" ($n = 81$; 23.2%).

**Assessment of therapists' attitudes toward routine outcome monitoring.** Therapists' attitudes toward routine outcome monitoring were assessed at the end of each treatment. An index of therapists' attitudes toward feedback was composed of two indicators: The first indicator was the number of modifications made due to feedback. After the treatments, therapists were asked what they actually did with the provided feedback. A list of 10 different options was presented to the therapists.[2] For each of these options therapists could choose whether they used the feedback in the respective way or not. Since the distribution of the reported number of modifications made due to feedback was skewed to the right (many made just "one modification" for a patient), patients were grouped by whether their therapists applied only one modification or more than one modification due to feedback. The second indicator was an item asking therapists' about their overall satisfaction with the project. Therapists' answered on a 5-point Likert scale ranging from 1 ("completely satisfied") to 5 ("dissatisfied"). The first three categories and the last two categories were pooled in a "satisfied" and a "not satisfied" group, respectively.

Subsequently, we combined these two indicators into one composite index. This index constitutes therapists' attitudes toward the feedback system in terms of frequency of usage and overall satisfaction with the system. Similar to the approach described above for the patients' attitudes index, separate categories for patients of the CG and those with missing values were built. Using the categories of this index the complete sample was divided into six groups: "satisfied/one modification" ($n = 57$; 16.3%), "not satisfied/one modification" ($n = 19$; 5.4%), "satisfied/several modifications" ($n = 107$; 30.7%), "not satisfied/several modifications" ($n = 26$; 7.4%), "missing" ($n = 59$; 16.9%), and "control group" ($n = 81$; 23.2%).

### Data Analysis Strategy

Since the data have a hierarchical structure with patients nested within therapists, multilevel models were used. Patients who are treated by the same therapist are likely to have more similar experiences than two randomly chosen patients treated by different therapists. Consequently, individual observations are not independent from each other, violating the assumption of independence. For these kinds of hierarchical data structure, multilevel modeling (MLM) has been established to be the method of choice because it addresses the fact that observations are not independent (e.g., Gallop & Tasca, 2014; Raudenbush & Bryk, 2002). For analyzing the present data, two-level models were used with patients on level 1 and therapists on level 2 (equations are reported in Appendix 1). The two-level model partitions the total variability into two components: variation within therapists (level 1) and variation between therapists (level 2). These components allow calculating the proportion of variance associated with each of the levels. The share of the total variance which is associated with level 2 is referred to in the literature as intraclass correlation or *therapist effect* (Baldwin & Imel, 2013). The larger the therapist effect the larger the differences between therapists concerning the dependent variable.

The analysis started in both cases with a *null model* that included no predictors. The other models were hierarchically built: Treatment outcome ($GSI_{post}$) was predicted by the following fixed effect predictors, each added in separate modeling steps on level one: initial symptom impairment ($GSI_{pre}$; *Model 1*), initial interpersonal problems ($IIP_{pre}$), early alliance ($HAQ_{pre}$), number of diagnoses (*Model 2*), early feedback (*Model 3*), therapists' attitudes toward feedback (*Model 4*), and patients' attitudes toward feedback (*Model 5*). Since the distribution of the $GSI_{post}$ scores suggests non-normality, all analyses concerning this variable have been conducted with a transformed (root squared) version of this variable.

Variance explained in the different steps is calculated by subtracting the residual variation of the respective model from the residual variation of an identical model not including the newly added predictor variables and dividing this difference through the residual variance of the model with fewer predictor variables (Raudenbush & Bryk, 2002).

The same approach was used to predict treatment length (log-transformed to the base $e$). The same predictors were used and added on level one, only the number of approved sessions by the insurance company was added as an additional control variable in Model 1. Given that the study arms were confounded with differences in the number of approved sessions, this variable is an important contextual design factor, with the potential to influence differences in the length of the treatments. Therefore, in the analyses on treatment length, we

controlled for the number of approved sessions when investigating the influence of other predictors. Data analyses were conducted with the free software environment R version 3.1.1 (R Development Core Team, 2014). The multilevel models have been estimated with the R package lme4 (Bates, Maechler, Bolker, & Walker, 2014).

## Results

### Comparison of the Full Sample and the Sample with at Least Five Patients per Therapist

To check for systematic selection effects the different subsamples were compared with regard to several variables. In a first step, the study patient sample ($N = 349$) was compared to the full patient sample ($N = 751$) and the sample of excluded patients ($N = 402$) regarding pre and post GSI scores, socio-demographic characteristics, as well as the predictor variables.

Figure 1 shows the pre- and post-distributions of GSI scores in density plots for those three samples. The three distributions are very similar at both time points. Additionally, Table I shows that the three samples do not differ significantly (overlapping 95% confidence intervals) with regard to socio-demographic variables (age and sex) as well as all the predictor and outcome variables included in the statistical analyses. The IG had significantly more approved and actually provided treatment sessions than the CG in each of the three samples. These differences resulted from the specific design of the study and were controlled in the following analyses.

Similarly, Table II shows that also the therapists from the three samples do not differ significantly (overlapping 95% confidence intervals) with regard to socio-demographic variables (age and sex), years of experience, proportion of cognitive behavioral therapists in the sample, working hours per week, average symptomatic intake impairment of their patients in the respective sample ($BSI_{pre}$), and the average interpersonal distress at intake of their patients in the respective sample ($IIP_{pre}$). Descriptively, the only difference between the samples is that more CBT therapists are in the study sample (71.8%) than in the excluded sample (59.8%).

### Prediction of Treatment Outcome

The results of the MLM are reported in Table III. *Model 1* is the widely used model for the quantification of therapist differences taking only into account initial patient impairment ($GSI_{pre}$; Baldwin & Imel, 2013; Saxon & Barkham, 2012). The initial GSI score explained 22.40% of the variance in patient-rated outcomes (compared to the *null model*). A patient with an average $GSI_{pre}$ score is on average 0.45 $GSI_{post}$ scores less impaired than a patient who starts with a one point higher $GSI_{pre}$ score. Dividing the residual level 2 variance (therapist variation) by the total variation results in a therapist effect of 5.88%. *Model 2 added* three additional intake characteristics: initial interpersonal distress ($IIP_{pre}$), early patient-rated alliance ($HAQ_{pre}$), and number of diagnoses[3] (*Model 2*). Only $HAQ_{pre}$ and the number of diagnoses were significant predictors of treatment outcome and were kept as predictors in the subsequent models. These predictors explain incrementally 6.53% of the variance on the patient level. In *Model 3* feedback (dummy coded as "negative feedback" and "no negative feedback" with the control group patients as reference category) was included. Feedback explained 10.62% of the variance in patient reported outcome. Receiving a negative feedback sign at session 10 resulted for an averagely impaired patient in 0.35 points higher $GSI_{post}$ scores compared to patients from the CG. Patients for whom their therapists receive no negative feedback over the course of treatment did not differ significantly from the CG with regard to their $GSI_{post}$ scores. Thus, in
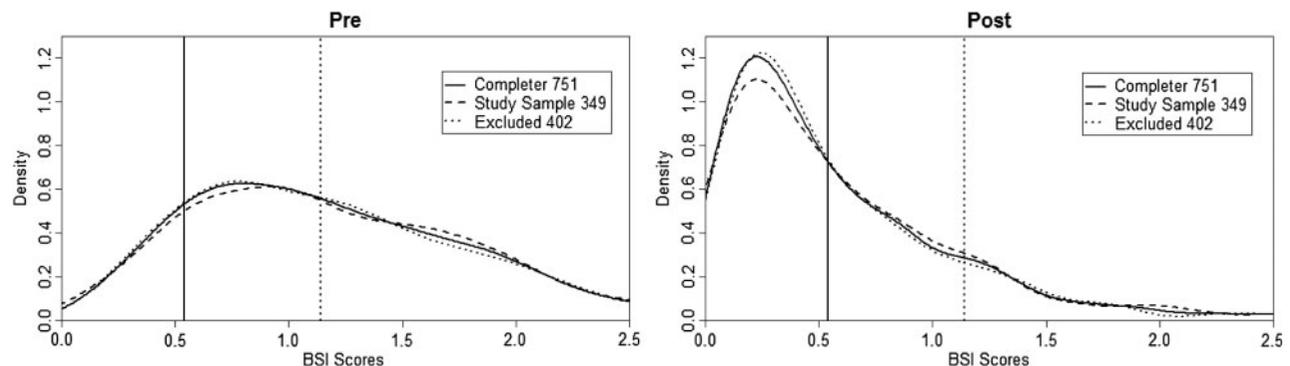


Figure 1. Distribution of BSI scores for the full sample ($N = 751$) and the subsample with at least 5 patients per therapist ($N = 349$) in relation to the cut-off between a non-clinical sample and an outpatient sample (solid vertical line) and the cut-off between an outpatient and an inpatient sample (dashed vertical line).

Table II. Descriptive therapist characteristics (Mean [*SD*]) and 95% CI for the completer sample, the study sample, and the sample of excluded patients for therapists grouped by treatment condition and on average.

| | Completer sample | | Study sample | | Excluded sample | |
|---|---|---|---|---|---|---|
| | IG (*n* = 90) | CG (*n* = 87) | IG (*n* = 31) | CG (*n* = 13) | IG (*n* = 59) | CG (*n* = 74) |
| Age | 48.5 (5.93) | 48.5 (5.23) | 49.8 (5.97) | 49.0 (3.94) | 47.9 (5.8) | 48.4 (5.4) |
| | 48.5 [47.64, 49.38] | | 49.6 [47.84, 51.36] | | 48.2 [47.15, 49.16] | |
| Female | 45.6% | 43.2% | 51.6% | 40.0% | 42.4% | 43.8% |
| | 44.5% [37.1%, 52.2%] | | 48.8% [34.3%, 63.5%] | | 43.1% [34.7%, 51.9%] | |
| Years of experience | 17.5 (6.20) | 16.5 (5.09) | 19.1 (6.37) | 17.5 (4.95) | 16.7 (6.0) | 16.4 (5.1) |
| | 17.1 [16.21, 17.98] | | 18.7 [16.80, 20.61] | | 16.6 [15.56, 17.54] | |
| Cognitive Behavior Therapy | 58.6% | 68.1% | 77.4% | 50.0% | 48.2% | 70.5% |
| | 62.8% [55.0%, 70.0%] | | 71.8% [56.1%, 83.6%] | | 59.8% [50.8%, 68.3%] | |
| Working hours per week | 37.9 (10.05) | 37.1 (9.59) | 40.9 (8.26) | 39.1 (8.95) | 36.4 (10.6) | 36.8 (9.7) |
| | 37.6 [36.05, 39.08] | | 40.4 [37.80, 43.08] | | 36.6 [34.80, 38.42] | |
| $BSI_{pre}$ | 1.24 (.44) | 1.13 (.49) | 1.26 (.31) | 1.11 (.29) | 1.23 (.49) | 1.13 (.52) |
| | 1.18 [1.11, 1.25] | | 1.22 [1.12, 1.31] | | 1.17 [1.09, 1.26] | |
| $IIP_{pre}$ | 1.56 (.39) | 1.47 (.44) | 1.54 (.23) | 1.51 (.22) | 1.58 (.46) | 1.46 (.47) |
| | 1.52 [1.45, 1.58] | | 1.53 [1.46, 1.60] | | 1.51 [1.43, 1.59] | |

*Note.* All values were calculated for available cases; missings were deleted. Mean scores are shown for age, years of experience, working hours per week, as well as for aggregated $BSI_{pre}$, and $IIP_{pre}$ scores; *SD*s in brackets. Female and CBT show percentages.

*Models 4 and 5* a combination of the CG and the "no negative feedback" group serves as reference group.

In the next step, therapists' attitudes toward feedback were investigated (*Model 4*). This variable explained additionally 5.39% of the variation in patients' outcomes. Especially the group of patients for whom therapists were satisfied with the project and made only one modification due to the feedback showed significantly lower $GSI_{post}$ scores (for an averagely impaired patient 0.17 points lower) than the control group.

In the last step, the influence of patients' attitudes toward feedback on treatment outcome was investigated (*Model 5*). This variable explained incrementally 5.72% of the variance in patients' outcomes. The group of patients that reported a positive attitude toward the feedback showed significantly lower $GSI_{post}$ scores than the control

Table III. MLM results for models predicting treatment outcome ($\sqrt{GSI_{post}}$).

| Models | Null model | Model 1 | Model 2 | Model 3 | Model 4 | Model 5 |
|---|---|---|---|---|---|---|
| *Fixed part* | | | | | | |
| Intercept | 0.70*** | 0.70*** | 0.70*** | 0.66*** | 0.68*** | 0.68*** |
| $GSI_{pre}$ | | 0.25*** | 0.22*** | 0.23*** | 0.23*** | 0.23*** |
| $IIP_{pre}$ | | | 0.02 | – | – | – |
| $HAQ_{pre}$ | | | −0.01* | −0.01* | −0.01** | −0.004+ |
| Number of diagnoses | | | 0.06** | 0.04+ | 0.03 | 0.04+ |
| Feedback | | | | | | |
| Negative feedback | | | | 0.20*** | 0.22*** | 0.24*** |
| No negative feedback | | | | −0.03 | – | – |
| Therapist attitude | | | | | | |
| Not satisfied/several mod. | | | | | 0.13+ | – |
| Not satisfied/1 mod. | | | | | −0.03 | – |
| Satisfied/several mod. | | | | | −0.05 | – |
| Satisfied/1 mod. | | | | | −0.12* | – |
| Missing | | | | | −0.08 | – |
| Patient attitude | | | | | | |
| Positive | | | | | | −0.25*** |
| Average/negative | | | | | | −0.08+ |
| Missing | | | | | | 0.02 |
| *Random part* | | | | | | |
| Level 2 | 0.009 | 0.005 | 0.007 | 0.006 | 0.007 | 0.002 |
| Level 1 | 0.103 | 0.080 | 0.075 | 0.067 | 0.063 | 0.063 |
| *Explained variance* | | | | | | |
| Level 1 (%) | | 22.4 | 27.5 | 35.1 | 38.8 | 39.0 |

*Note.* N = 349 patients, 44 therapists; Feedback neutral is not significantly different from the control group and serves together with the control group in further analysis as reference category.
***$p$ = 0; **$p$ < 0.001; *$p$ < 0.05; +$p$ < 0.1.

group (for an averagely impaired patient 0.32 points).

### Relation of Treatment Outcome and Attitudes toward Feedback

To explore the relation between outcome and attitudes in more detail and facilitate interpretation, Figures 2 and 3 show the adjusted pre--post effect sizes (adjusted *d*) for the different therapist and patient attitude groups, respectively. Separate MLM with standardized (at the mean and SD of the $GSI_{pre}$ scores) $GSI_{post}$ scores as dependent
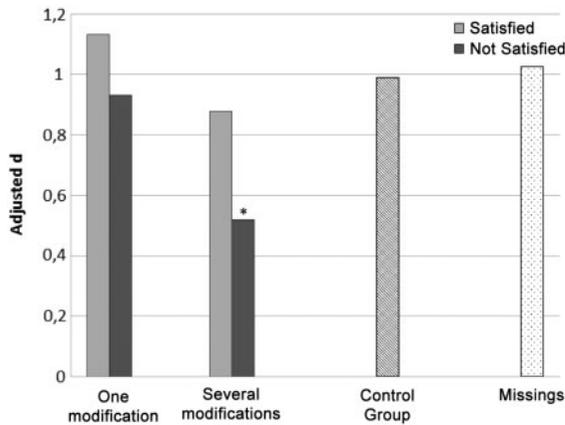


Figure 2. Effect sizes (*d*) for the different groups of therapists' attitudes toward feedback, the control group, and those patients for whom no information is reported (missings) adjusted for initial impairment ($GSI_{pre}$), early alliance ($HAQ_{pre}$), and number of diagnoses (*$\star p$* < 0.05).

variable were conducted to produce adjusted pre–post effect sizes and to compare the different levels of therapists' and patients' attitudes. Due to the standardization of the variables, the adjusted means can be interpreted as adjusted pre–post effect sizes (adjusted *d*). Figure 2 depicts the adjusted pre–post effects for the different therapist attitude groups, the control group, and the missing group. *MLM* revealed that the effects sizes of the group of patients for whom their therapists made several modifications due to feedback and are not satisfied with the monitoring system ("not satisfied/several modifications") are significantly lower than for all other attitude groups. Patients in this group had an effect sizes of $d_{adjusted}$ = .52. Patients for whom their therapists were satisfied with the project and made only one modification due to feedback ("one modification/satisfied") had the highest adjusted effect size ($d_{adjusted}$ = 1.13). Patients for whom their therapists were satisfied with the project and made several modifications due to feedback ("several modifications/satisfied") had the second to lowest adjusted effect size ($d_{adjusted}$ = 0.88) which was similar to the adjusted effect size of patients for whom their therapist were not satisfied and made only one modification due to feedback ($d_{adjusted}$ = 0.93). Patients from the CG ($d_{adjusted}$ = 0.99) had a similar adjusted effect size as patients with missing values on these variables ($d_{adjusted}$ = 1.00).

Figure 3 depicts the adjusted means of the pre–post effects sizes (adjusted *d*) for the different patient attitude groups, the control group, and the missing
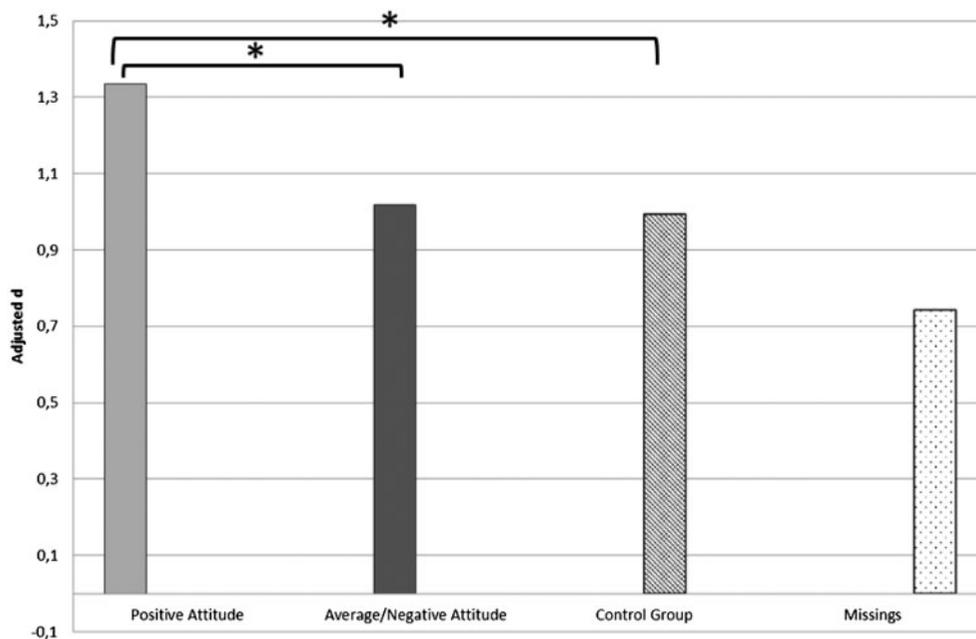


Figure 3. Effect sizes (*d*) for the different groups of patients' attitudes toward feedback, the control group, and those patients for whom no information is reported (missings) adjusted for initial impairment ($GSI_{pre}$), early alliance ($HAQ_{pre}$), and number of diagnoses (*$\star p$* < 0.05).

group. Patients who reported positive attitudes toward feedback had descriptively the highest adjusted effect size ($d_{adjusted}$ = 1.33). *MLM* revealed that the positive attitude group was significantly different from the average/negative attitude group ($d_{adjusted}$ = 1.02) and the control group ($d_{adjusted}$ = 0.99). The missing group had significantly lower adjusted effect sizes than all other groups ($d_{adjusted}$ = 0.74).

## Prediction of Treatment Length

The approach for the prediction of treatment length (number of sessions) was similar to the one above[4] and five models with an increasing number of predictors were used to test the incremental impact of these predictors (see Table IV). *Model 1* included the number of approved sessions by the insurance company as single predictor of treatment length. About 20.03% of the differences in treatment length could be explained by the number of approved sessions. After controlling for approved sessions, we found a therapist effect of about 8.89%. In *Model 2* additional intake characteristics were included in the analysis. Together, these variables explained 2.09% incremental variance compared to *Model 1*. Only patients' mean scores in the IIP at the beginning of the treatment and *number of diagnoses* were

significant and marginally significant predictors, respectively. Thus, initial symptom impairment ($GSI_{pre}$) and early alliance ($HAQ_{pre}$) were excluded from subsequent analyses.

*Model 3* tested the additional influence of the feedback on treatment length. Feedback explained additional 1.3% of the variation in treatment length. Patients with negative feedback had significantly longer treatments than all other patients.

*Models 4* and *5* tested the additional influence of therapists' and patients' attitudes toward feedback on treatment length, respectively. Neither therapists' nor patients' attitudes toward feedback significantly explained additional variance in treatment length.

## Associations between Therapist Effects on Outcome and Length

Finally, the association between therapist effects on outcome and treatment length was investigated. Figure 4 shows the association between therapist residuals from the outcome and the treatment length prediction model. On the $x$ axis therapists' deviations from the average $\sqrt{GSI_{post}}$ score controlled for initial impairment (at the position of the average score, the residuals equal 0) are depicted. On the $y$ axis the deviations from the average treatment length controlled for approved sessions

Table IV. MLM results for models predicting treatment length (ln(number of sessions)).

| Models | Null model | Model 1 | Model 2 | Model 3 | Model 4 | Model 5 |
|---|---|---|---|---|---|---|
| *Fixed part* | | | | | | |
| Intercept | 3.55*** | 3.55*** | 3.56*** | 3.52*** | 3.58*** | 3.58*** |
| Approved session | | 0.02*** | 0.02*** | 0.02*** | 0.02*** | 0.02*** |
| $GSI_{pre}$ | | | −0.03 | – | – | – |
| $IIP_{pre}$ | | | 0.15* | 0.13** | 0.13** | 0.13** |
| $HAQ_{pre}$ | | | −0.00 | – | – | – |
| Number of diagnoses | | | 0.07+ | 0.05 | 0.06 | 0.06 |
| Feedback | | | | | | |
| Negative feedback | | | | 0.14* | 0.15* | 0.15* |
| Therapist attitude | | | | | | |
| Not satisfied/several mod. | | | | | 0.04 | – |
| Not satisfied/1 mod. | | | | | −0.04 | – |
| Satisfied/several mod. | | | | | −0.10 | – |
| Satisfied/1 mod. | | | | | −0.21* | – |
| Missing | | | | | −0.06 | – |
| Patient attitude | | | | | | |
| Positive | | | | | | −0.07 |
| Average/negative | | | | | | −0.02 |
| Missing | | | | | | −0.15+ |
| *Random part* | | | | | | |
| Level 2 | 0.040 | 0.020 | 0.015 | 0.015 | 0.010 | 0.014 |
| Level 1 | 0.256 | 0.205 | 0.200 | 0.197 | 0.199 | 0.194 |
| *Explained variance* | | | | | | |
| Level 1 | | 20.03 | 21.7 | 23.0 | 22.7 | 24.3 |

*Note.* N = 336 patients, 44 therapists; Feedback neutral is not significantly different from the control group and serves together with the control group in further analysis as reference category.
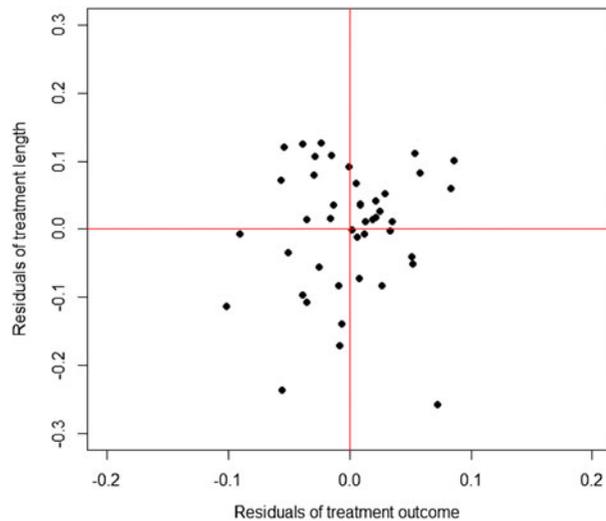***$p$ = 0; **$p$ < 0.001; *$p$ < 0.05; +$p$ < 0.1.

Figure 4. Scatterplot of therapist intercept residuals for GSI outcome scores and treatment length produced by *Model 1* for outcome and length, respectively. These represent how much each therapist's outcome and treatment length across their patients differ from the average treatment outcome and treatment length in this data-set, controlling for initial impairment as well as approved sessions. Negative residuals represent better outcomes as well as shorter treatments.

(at the position of the average score, the residuals equal 0) are depicted. Each point in the scatterplot represents the average $\sqrt{GSI_{post}}$ score of a therapist in reference to the average $\sqrt{GSI_{post}}$ score from all therapists as well as the average treatment length of a therapist in reference to the average treatment length from all therapists in our analysis sample. Thus, the intersection of the axes marks the concurrence of average effectivity (which relates to a score of $GSI_{post} = 0.49$) and average treatment length (which relates to an average treatment length of 35.0 sessions) per therapist. Therapists located in the upper left quadrant of the scatterplot are those who have on average relatively positive outcomes and rather long treatments. Therapists with good outcomes and short treatments are located in the bottom left quadrant, those with relatively bad outcomes and long treatments in the upper right quadrant, and those with bad outcomes and short treatments in the bottom right quadrant. There was no significant correlation between the average length of the treatments and the effectiveness of a therapist ($r = .09$; $p = 0.56$).

## Discussion

Recent investigations which examined the effects of routine outcome monitoring and feedback suggest that feedback does not work uniformly for every patient and every therapist (De Jong et al., 2012; Simon et al., 2012). The present study investigated therapist effects on and predictors of treatment outcome and treatment length in a naturalistic setting. Therapist differences were found to account for 6.2% of the differences in patients' outcomes and for 9% of the differences in treatment lengths. There was no significant relation between therapists' average effectiveness and their average treatment length. Therapists' and patients' attitudes toward feedback both explained an additional share of 5.4% (therapist attitudes) and 5.7% (patient attitudes) of the variation in outcomes, respectively. For patients, a positive attitude was associated with better outcomes. For therapists it showed to be most promising when they were satisfied with the system and used the feedback for one specific modification per patient. Neither therapists' nor patients' attitudes toward feedback could incrementally explain significant variation in treatment length.

Furthermore, higher initial symptomatic impairment, lower early patient-rated alliance, and several diagnoses were consistently associated with worse outcome while initial interpersonal distress showed no significant relation. In accordance with previous studies negative feedback in the IG at session 10 was a significant predictor of worse treatment outcome. Although not on-track patients show improvements later on, they do not reach the same level at the end of therapy as on-track patients. A little less variation could be explained in treatment length compared to treatment outcomes (about 24% vs. 39%). Differences in treatment length were predicted by the number of approved sessions, by different intake levels of interpersonal distress (more interpersonal distress, longer treatments), and by feedback (negative feedback, longer treatments).

The finding that deteriorating patients (negative feedback at session 10) have worse treatment outcomes is well known. Previous studies showed that feedback is especially helpful for deteriorating patients even so not-on track patients do not reach the same outcome level as on-track patients (e.g., Lambert, 2007). Our study differs from previous efforts insofar as usually deteriorating patients from an IG (continuous assessments plus feedback to therapists) were compared with deteriorating patients from a control group (continuous assessments without feedback to therapists). In contrast to that, the current investigation compared deteriorating patients from the IG (continuous assessments plus feedback to therapists) with all patients from the CG for which only pre- and post-assessments were available. Without continuous assessments over the course of treatment in the CG, we could not match deteriorating patients from the CG to those from the IG (Strauss et al., 2015). Despite the different kinds of CGs, our findings are not in contrast to the

findings of previous studies and our results emphasize the general importance to track treatment outcome continuously over the course of treatment and refer to the importance to provide feedback especially to those patients who show negative developments. Also in our study a more positive development after psychometric feedback was detected for the not-on track patients in the IG (discussed in more detail in Lutz et al., 2013).

Therapists' and patients' attitudes toward feedback showed to be potentially important variables to consider in future feedback studies. Regarding therapists' attitudes, especially those patients for whom their therapists were satisfied with the feedback system and who made only one modification due to the feedback showed the best outcomes. Patients showed worse outcomes when therapists were not satisfied with the feedback system and yet made several modifications due to the feedback. However, it should be noted that these associations are just correlational and not causal, since modifications and satisfaction were assessed at the end of each treatment. One specific modification (after feedback) seems to contribute to a positive outcome, if therapists are satisfied with the feedback. However, with regard to the post-hoc assessment of therapist satisfaction with the feedback system an alternative explanation might be that the outcome predicted the attitudes. Therapist and patients who had a good experience with therapy might look back on the feedback positively, especially those therapists who only had to make one modification, while those who had a worse experience view the feedback more negatively. Besides the impact of therapists' attitudes, results also suggest an influence of patients' attitudes toward feedback. It could be an important factor for the success of a treatment to provide patients with a sensible rationale for why these continuous assessments are an integral and important part of clinical practice (see also Boswell, Kraus, Miller, & Lambert, 2013; Castonguay et al., 2013; Flückiger et al., 2013).

The fact that therapists differ in both their attitudes as well as their use of feedback raises questions of treatment integrity within feedback studies (Wittmann & Lutz, 2014). Similarly to adhering to a treatment manual it should be checked whether therapists understand and actually integrate feedback information in their treatments when participating in the intervention group of such a study. Therefore, it is important to train therapists in the use of feedback instruments, document what they actually did with the feedback, and assess their attitude toward routine outcome monitoring.

Regarding treatment length it came as a surprise to us that higher initial symptom distress was not associated with more treatment sessions, when controlled for approved sessions. Therefore, we tested if this was a mere consequence of the fact that generally more sessions were approved for more distressed patients. If that would be the case the effect of initial symptom distress on treatment length would have been superimposed by the differential number of approved sessions. However, no significant relation could be found between $GSI_{pre}$ scores and the number of approved sessions.

Another central aim of the present study was to quantify the therapist effects on treatment outcome and length. The variance on patient-rated outcome that could be attributed to differences between therapists (6.2%) was in the range of previous studies and meta-analysis in this field (Baldwin & Imel, 2013). About 9% of patients' differences in treatment length are due to differences between therapists. Since this is the first study investigating therapist effects on treatment lengths it is difficult to evaluate the absolute size of this effect. Nevertheless, it seems safe to say that more research in naturalistic samples on this currently neglected aspect could be worthwhile. Besides raising the question of why therapists differ in this regard (even after controlling for differences in their patients initial symptom distress), results from such studies could elucidate cognitive models and other determinants of therapists' choices. An incorporation of these models into more regulated treatment settings (e.g., clinical trials or structured intervention programs) could make these programs more amenable and possibly increase the acceptance on the side of the therapists.

It would be important for future research to replicate this finding in other samples, settings and potentially with other modeling strategies. Especially discrete-time event hazard models (e.g., Hox, 2010) might be an interesting alternative for the analysis of treatment length. We chose the presented modeling strategy in this study to build the treatment length models as similar to the established treatment outcome models as possible.

No relation between therapists' average effectiveness and their average treatment length was found. This suggests that the most effective therapists are not those who on average provide shorter or longer treatments. Although therapists' effectiveness and treatment lengths were unrelated, a joined investigation of therapist effects on outcome and length allows the identification of specific relations on a therapists' level (see Figure 4). The use of multilevel models also allows to control for the individual therapist's (or service's) case-mix, which can – as pointed out above – have an influence on patient outcome as well as treatment length.

A next step could be a closer look at the most effective therapists by means of video analyses, as well as patient and therapist interviews. There is a huge gap between our knowledge of how influential therapists are with regard to outcome and length of treatments and our knowledge of what makes therapists that influential. The aim might be to extract those elements that make therapists effective and transform them in parts of established clinical training programs (Castonguay, Eubanks, Goldfried, Muran, & Lutz, 2015). For health service research it is important to note that more effective therapists do not seem to provide generally shorter or longer treatments. Successful therapists seem to adapt their treatment length to the specific patient and the specific progress pattern. However, given that our findings are based on a limited number of therapists (*N* = 44), the results from these analyses should be treated with caution.

In addition to the limitations noted above, some further aspects limit the scope of this study and give suggestions for future research. Due to resistance among therapists and the study's steering committee, it was not possible to implement the standard control group of these studies in which patients filled out psychometric questionnaires over the course of the treatment but no feedback was given (Strauss et al., 2015). This CG would have allowed comparisons with the IG within patient subgroups showing similar change courses (e.g., deteriorating patients). Future research implementing such a design can look at differences in therapists' and patients' attitudes toward routine outcome monitoring among groups with similar change courses (on track/not on track).

Another issue is connected to the post-hoc assessment of the therapists' attitudes toward feedback. Since the feedback procedure was provided before the therapists' attitudes were assessed it is possible that the kind of feedback the therapists received over the course of the treatment influenced their attitudes. However, additional analyses for the IG revealed that the proportion of positive, neutral, or negative feedback (taking into account the total number of feedback reports given for this patient) had no statistically significant influence on either therapists' satisfaction or the number of modifications made due to feedback.

Furthermore, it should be noted that the diagnostic procedure was different between the CG and the IG. While in the CG diagnosis was based on clinical judgment, in the IG a standardized assessment procedure was conducted including a structured diagnostic interview (related to ICD-10 criteria; Hiller et al., 2004). This difference led to a slightly higher average number of diagnoses for the patients in the IG (see Table I), therefore we included the number of diagnoses as a potential covariate in all our prediction models.

Another shortcoming of the presented analyses is that the sample inclusion criteria were rather strict. It could be that those patients who met the inclusion criteria of this study (completion of intake and post-treatment assessment; at least 10 sessions) were rather satisfied with the treatment or are in other regards a specific, non-representative sample. However, we showed in this and previous reports that the study sample is similar to the completer sample and that the completer sample is similar to patients undergoing outpatient psychotherapy in Germany (Lutz, Wittmann, et al., 2013). Nevertheless, a replication could shed further light on the robustness of our findings.

Despite these limitations, the findings of this study add to the still growing body of literature on feedback and therapist effects. It seems that for the application of feedback tools in routine care it is necessary not only to train therapists in how to use the feedback, but also to convince them of the value these tools add by supporting (not replacing) their everyday clinical decision-making.

## Disclosure statement

No potential conflict of interest was reported by the authors.

## Notes

[1] Since not enough data were available for the psychoanalytic treatment modality, the present study includes only CBT and psychodynamic treatments.

[2] This list was composed of the statement "Due to the feedback, I …" and the following 10 different options from which multiple choices were possible: "… discussed with the patient his/her answers in the questionnaire"; "… tried to assist the patient's resources"; "… tried to adjust my therapeutic interventions"; "… discussed with the patient his/her interpersonal problems"; "… prepared the end of the therapy"; "… tried to enhance the patient's motivation for therapy"; "… varied the intervals between sessions; "… tried to enhance the therapeutic alliance"; "… consulted additional sources of help (e.g., supervision, intervision, literature, further education); "… tried new homework with the patient."

[3] In a preliminary analysis, we checked if differences in diagnoses had an impact on outcome. In this analysis, seven dummy variables were created with depression as reference group: dysthymia, panic disorder, other anxiety disorders, adjustment disorder, eating disorder, personality disorders, and other disorders. Adding the dummy variables as predictors did not explain significant variation in $\sqrt{\text{GSI}_{\text{post}}}$ scores. Since the diagnostic category was not a significant predictor of treatment outcome and too many dummy variables in one model led to instable model estimations, we did not leave these non-significant predictors in next modeling step.

[4] For 13 patients from the analysis sample (*N* = 349) no information regarding treatment length was given and therefore had to be excluded (*N* = 336).

# References

Alexander, L. B., & Luborsky, L. (1986). The Penn Helping Alliance Scales. In L. S. Greenberg & W. M. Pinsof (Eds.), *The psychotherapeutic process: A research handbook* (pp. 325–366). New York, NY: Guilford Press.

Baldwin, S. A., & Imel, Z. E. (2013). Therapist effects: Findings and methods. In M. J. Lambert (Ed.), *Bergin and Garfield's handbook of psychotherapy and behavior change* (6th ed., pp. 85–133). New York, NY: John Wiley & Sons.

Baldwin, S. A., Murray, D. M., Shadish, W. R., Pals, S. L., Holland, J. M., Abramowitz, J. S., … Watson, J. (2011). Intraclass correlation associated with therapists: Estimates and applications in planning psychotherapy research. *Cognitive Behaviour Therapy, 40*(1), 15–33. doi:10.1080/16506073.2010.520731

Bassler, M., Potratz, B., & Krauthauser, H. (1995). Der "Helping Alliance Questionnaire" (HAQ) von Luborsky. Möglichkeiten zur Evaluation des therapeutischen Prozesses von stationärer Psychotherapie [The 'Helping Alliance Questionnaire' (HAQ) by Luborsky]. *Psychotherapeut, 40*, 23–32.

Bates, D., Maechler, M., Bolker, B., & Walker, S. (2014). lme4: Linear mixed-effects models using Eigen and S4 (R package version 1.1–7). Retrieved from http://CRAN.R-project.org/package=lme4

Beck, A. T., Steer, R. A., & Carbin, M. G. (1988). Psychometric properties of the Beck Depression Inventory: Twenty-five years of evaluation. *Clinical Psychology Review, 8*(1), 77–100. doi:10.1016/0272-7358(88)90050-5

Boswell, J. F., Kraus, D. R., Miller, S. D., & Lambert, M. J. (2013). Implementing routine outcome monitoring in clinical practice: Benefits, challenges, and solutions. *Psychotherapy Research, 25*, 6–19. doi:10.1080/10503307.2013.817696

Byrne, S. L., Hooke, G. R., Newnham, E. A., & Page, A. C. (2012). The effects of progress monitoring on subsequent readmission to psychiatric care: A six-month follow-up. *Journal of Affective Disorders, 137*(1–3), 113–116. doi:10.1016/j.jad.2011.12.005

Carlier, I. V. E., Meuldijk, D., Van Vliet, I. M., Van Fenema, E., Van der Wee, N. J. A., & Zitman, F. G. (2012). Routine outcome monitoring and feedback on physical or mental health status: Evidence and theory. *Journal of Evaluation in Clinical Practice, 18*(1), 104–110. doi:10.1111/j.1365-2753.2010.01543.x

Castonguay, L. C., Barkham, M., Lutz, W., & McAleavy, A. (2013). Practice-oriented research: Approaches and applications. In M. J. Lambert (Ed.), *Bergin and Garfield's handbook of psychotherapy and behavior change* (6th ed., pp. 85–133). New York, NY: John Wiley & Sons.

Castonguay, L. G., Eubanks, C. F., Goldfried, M. R., Muran, J. C., & Lutz, W. (2015). Research on psychotherapy integration: Building on the past, looking to the future. *Psychotherapy Research, 25*, 365–382.

Crits-Christoph, P., Baranackie, K., Kurcias, J., Beck, A., Carroll, K., Perry, K., … Zitrin, C. (1991). Meta-analysis of therapist effects in psychotherapy outcome studies. *Psychotherapy Research, 1*(2), 81–91. doi:10.1080/10503309112331335511

De Jong, K., Timman, R., Hakkaart-Van Roijen, L., Vermeulen, P., Kooiman, K., Passchier, J., & Busschbach, J. V. (2014). The effect of outcome monitoring feedback to clinicians and patients in short and long-term psychotherapy: A randomized controlled trial. *Psychotherapy Research, 24*, 1–11.

De Jong, K., van Sluis, P., Nugter, M. A., Heiser, W. J., & Spinhoven, P. (2012). Understanding the differential impact of outcome monitoring: Therapist variables that moderate feedback effects in a randomized clinical trial. *Psychotherapy Research, 22*, 464–474. doi:10.1080/10503307.2012.673023

Derogatis, L. R. (1975). *Brief symptom inventory*. Baltimore, MD: Clinical Psychometric Research.

Ehlers, A., & Margraf, J. (2001). *Fragebogen zu körperbezogenen Ängsten, Kognitionen und Vermeidung (AKV)* [Questionnaire on body-directed anxieties, cognitions, and avoidance] (2., über-arbeitete und neunormierte Auflage). Weinheim: Beltz.

Flückiger, C., Holforth, M. G., Znoj, H. J., Caspar, F., & Wampold, B. E. (2013). Is the relation between early post-session reports and treatment outcome an epiphenomenon of intake distress and early response? A multi-predictor analysis in outpatient psychotherapy. *Psychotherapy Research, 23*(1), 1–13.

Franke, G. (2000). *BSI. Brief Symptom Inventory: Deutsche Version (BSI. Brief Symptom Inventory: German Version). Manual*. Göttingen: Beltz.

Gallop, R., & Tasca, G. A. (2014). Multilevel modeling of longitudinal data for psychotherapy researchers: 2. The complexities. In W. Lutz & S. Knox (Eds.), *Quantitative and qualitative methods in psychotherapy research* (pp. 117–141). New York, NY: Routledge.

Garner, D. M. (1991). *Eating disorders inventory-2: Professional manual*. Lutz, FL: Psychological Assessment Resources.

Hannan, C., Lambert, M. J., Harmon, C., Nielsen, S. L., Smart, D. W., Shimokawa, K., & Sutton, S. W. (2005). A lab test and algorithms for identifying clients at risk for treatment failure. *Journal of Clinical Psychology, 61*, 155–163. doi:10.1002/jclp.20108

Hiller, W., Zaudig, M., & Mombour, W. (2004). *IDCL – International Diagnostic Checklists for ICD-10 and DSM-IV*. Bern: Huber.

Horowitz, L. M., Strauß, B., & Kordy, H. (2000). *Inventar Interpersonaler Probleme (IIP-D): Deutsche Version [Inventory of interpersonal problems (IIP-D): German Version]*. Weinheim: Beltz Test GmbH.

Hox, J. (2010). *Multilevel analysis: Techniques and applications*. New York: Routledge.

Howard, K. I., Moras, K., Brill, P. L., Martinovich, Z., & Lutz, W. (1996). Evaluation of psychotherapy. Efficacy, effectiveness, and patient progress. *The American psychologist, 51*, 1059–1064. doi:10.1037/0003-066X.51.10.1059

Jacobson, N. S., & Truax, P. (1991). Clinical significance: A statistical approach to defining meaningful change in psychotherapy research. *Journal of Consulting and Clinical Psychology, 59*(1), 12–19. doi:10.1037/0022-006X.59.1.12

Klepsch, R., Zaworka, W., Hand, I., Lünenschloss, K., & Jauering, G. (1993). *Hamburger Zwangsinventar—Kurzform (HZI-K). Manual* [Hamburg Inventory for obsessive-compulsive disorders—short form]. *Göttingen: Beltz*.

Lambert, M. J. (2007). Presidential address: What we have learned from a decade of research aimed at improving psychotherapy outcome in routine care. *Psychotherapy Research, 17*(1), 1–14. doi:10.1080/10503300601032506

Lambert, M. J., & Shimokawa, K. (2011). Collecting client feedback. *Psychotherapy, 48*(1), 72–79. doi:10.1037/a0022238

Lambert, M. J., Whipple, J. L., Hawkins, E. J., Vermeersch, D. A., Nielsen, S. L., & Smart, D. W. (2003). Is it time for clinicians to routinely track patient outcome? A meta-analysis. *Clinical Psychology: Science and Practice, 10*, 288–301. doi:10.1093/clipsy.bpg025

Lutz, W., Böhnke, J. R., & Köck, K. (2011). Lending an ear to feedback systems: Evaluation of recovery and non-response in psychotherapy in a German outpatient setting. *Community Mental Health Journal, 47*, 311–317. doi:10.1007/s10597-010-9307-3

Lutz, W., Ehrlich, T., Rubel, J., Hallwachs, N., Röttger, M.-A., Jorasz, C., … Tschitsaz-Stucki, A. (2013). The ups and downs of psychotherapy: Sudden gains and sudden losses identified

with session reports. *Psychotherapy Research, 23*(1), 14–24. doi:10.1080/10503307.2012.693837

Lutz, W., Hofmann, S. G., Rubel, J., Boswell, J. F., Shear, M. K., Gorman, J. M., … Barlow, D. H. (2014). Patterns of early change and their relationship to outcome and early treatment termination in patients with panic disorder. *Journal of Consulting and Clinical Psychology, 82*, 287–297. doi:10.1037/a0035535

Lutz, W., Leon, S. C., Martinovich, Z., Lyons, J. S., & Stiles, W. B. (2007). Therapist effects in outpatient psychotherapy: A three-level growth curve approach. *Journal of Counseling Psychology, 54*(1), 32–39. doi:10.1037/0022-0167.54.1.32

Lutz, W., Wittmann, W., Böhnke, J., Rubel, J., & Steffanowski, A. (2013). Zu den Ergebnissen des Modellprojektes der Techniker-Krankenkasse – Ein Plädoyer für mehr Psychotherapieforschung in Deutschland [The TK-project quality monitoring in outpatient psychotherapy from the perspective of the evaluation team – A plea for more psychotherapy research in Germany]. *PPmP-Psychotherapie Psychosomatik Medizinische Psychologie, 63*, 225–228.

Meehl, P. E. (1954). *Clinical versus statistical prediction: A theoretical analysis and a review of the evidence*. Minneapolis: University of Minnesota Press.

Newnham, E. A., & Page, A. C. (2010). Bridging the gap between best evidence and best practice in mental health. *Clinical Psychology Review, 30*(1), 127–142. doi:10.1016/j.cpr.2009.10.004

Okiishi, J., Lambert, M. J., Nielsen, S. L., & Ogles, B. M. (2003). Waiting for supershrink: An empirical analysis of therapist effects. *Clinical Psychology & Psychotherapy, 10*, 361–373. doi:10.1002/cpp.383

Poston, J. M., & Hanson, W. E. (2010). Meta-analysis of psychological assessment as a therapeutic intervention. *Psychological Assessment, 22*(2), 203.

Probst, T., Lambert, M. J., Loew, T. H., Dahlbender, R. W., Göllner, R., & Tritt, K. (2013). Feedback on patient progress and clinical support tools for therapists: Improved outcome for patients at risk of treatment failure in psychosomatic in-patient therapy under the conditions of routine practice. *Journal of Psychosomatic Research, 75*, 255–261. doi:10.1016/j.jpsychores.2013.07.003

R Development Core Team. (2014). *R: A language and environment for statistical computing*. Vienna: R Foundation for Statistical Computing. Retrieved from http://www.R-project.org/

Raudenbush, S., & Bryk, A. S. (2002). *Hierarchical linears models: Applications and data analysis methods* (2nd ed.). Newbury Park, CA: Sage.

Rief, W., Hiller, W., & Heuser, J. (1997). *SOMS—Das Screening für somatoforme Störungen (Manual zum Fragebogen)* [Screening for somatoform disorders]. *Bern: Huber.*

Riemer, M., & Bickman, L. (2011). Using program theory to link social psychology and program evaluation. In M. M. Mark, S. I. Donaldson, & B. Campbell (Eds.), *Social psychology and evaluation* (pp. 104–139). New York, NY: Guilford.

Saxon, D., & Barkham, M. (2012). Patterns of therapist variability: Therapist effects and the contribution of patient severity and risk. *Journal of Consulting and Clinical Psychology, 80*, 535–546. doi:10.1037/a0028898

Shimokawa, K., Lambert, M. J., & Smart, D. W. (2010). Enhancing treatment outcome of patients at risk of treatment failure: Meta-analytic and mega-analytic review of a psychotherapy quality assurance system. *Journal of Consulting and Clinical Psychology, 78*, 298–311. doi:10.1037/a0019247

Simon, W., Lambert, M. J., Busath, G., Vazquez, A., Berkeljon, A., Hyer, K., … Berrett, M. (2013). Effects of providing patient progress feedback and clinical support tools to psychotherapists in an inpatient eating disorders treatment program: A randomized controlled study. *Psychotherapy Research, 23*, 287–300. doi:10.1080/10503307.2013.787497

Simon, W., Lambert, M. J., Harris, M. W., Busath, G., & Vazquez, A. (2012). Providing patient progress information and clinical support tools to therapists: Effects on patients at risk of treatment failure. *Psychotherapy Research, 22*, 638–647. doi:10.1080/10503307.2012.698918

Strauss, B. M., Lutz, W., Steffanowski, A., Wittmann, W. W., Boehnke, J. R., Rubel, J., … Kirchmann, H. (2015). Benefits and challenges in practice-oriented psychotherapy research in Germany: The TK and the QS-PSY-BAY projects of quality assurance in outpatient psychotherapy. *Psychotherapy Research, 25*(1), 32–51. doi:10.1080/10503307.2013.856046

Stulz, N., Lutz, W., Leach, C., Lucock, M., & Barkham, M. (2007). Shapes of early change in psychotherapy under routine outpatient conditions. *Journal of Consulting and Clinical Psychology, 75*, 864–874. doi:10.1037/0022-006X.75.6.864

Wittmann, W. W., Lutz, W., Steffanowski, A., Kriz, D., Glahn, E. M., Völkle, M.C., … Ruprecht, T. (2011). *Qualitätsmonitoring in der ambulanten Psychotherapie: Modellprojekt der Techniker Krankenkasse – Abschlussbericht.* Hamburg: Techniker Krankenkasse.

Wittmann, W. W., & Lutz, W. (2014, June). *The experience of a computer-based feedback system from the clients` perspective*. Paper presented at the 45th International Meeting of the Society for Psychotherapy Research (SPR), Copenhagen, DK.

## Appendix 1

**Model 1 for predicting treatment outcome**

Level 1 (Patient level): $\sqrt{\text{GSI}_{\text{post }ij}} = \beta_{0j} + \beta_{1j} \times$ initial impairment_gm$_{ij}$ + $e_{ij}$

Level 2 (Therapist level): $\beta_{0j} = \gamma_{00} + r_{0j}$; $\beta_{1j} = \gamma_{10}$

*Note.* MLM *formulas* for Model 1 predicting treatment outcome (GSI) where patient $i$ is nested within therapist $j$. Initial impairment was included as predictor on level 1. For Models 2–5 additional variables were tested as predictors at level 1: GSI$_{\text{pre\_gm}}$, IIP$_{\text{pre\_gm}}$, HAQ$_{\text{pre\_gm}}$, number of diagnoses_gm, feedback, patient attitude, and therapist attitude. Variables with the suffix "_gm" are included as grand-mean centered.

**Model 1 for predicting treatment length**

Level 1 (Patient level): ln_treatment_length$_{ij}$ = $\beta_{0j} + \beta_{1j} \times$ approved sessions_gm$_{ij}$ + $e_{ij}$

Level 2 (Therapist level): $\beta_{0j} = \gamma_{00} + r_{0j}$; $\beta_{1j} = \gamma_{10}$

*Note.* MLM formulas for Model 1 predicting treatment length where patient $i$ is nested within therapist $j$. Approved sessions were included as predictor at level 1. For Models 2–5 additional variables were tested as predictors at level 1: GSI$_{\text{pre\_gm}}$, IIP$_{\text{pre\_gm}}$, HAQ$_{\text{pre\_gm}}$, number of diagnoses_gm, feedback, patient attitude, and therapist attitude.