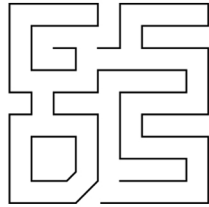
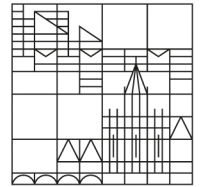


GRADUATE SCHOOL
OF DECISION SCIENCES



Universität
Konstanz



GSDS
Working Paper
No. 2017-05

Tracing the path of forgetting in rule abstraction and exemplar retrieval

Janina A. Hoffmann
Bettina von Helversen
Regina A. Weilbacher
Jörg Rieskamp

February 2017

Graduate School of Decision Sciences

All processes within our society are based on decisions – whether they are individual or collective decisions. Understanding how these decisions are made will provide the tools with which we can address the root causes of social science issues.

The GSDS offers an open and communicative academic environment for doctoral researchers who deal with issues of decision making and their application to important social science problems. It combines the perspectives of the various social science disciplines for a comprehensive understanding of human decision behavior and its economic and political consequences.

The GSDS primarily focuses on economics, political science and psychology, but also encompasses the complementary disciplines computer science, sociology and statistics. The GSDS is structured around four interdisciplinary research areas: (A) Behavioural Decision Making, (B) Intertemporal Choice and Markets, (C) Political Decisions and Institutions and (D) Information Processing and Statistical Analysis.

GSDS – Graduate School of Decision Sciences
University of Konstanz
Box 146
78457 Konstanz

Phone: +49 (0)7531 88 3633

Fax: +49 (0)7531 88 5193

E-mail: gds.office@uni-konstanz.de

-gds.uni-konstanz.de

ISSN: 2365-4120

February 2017

© 2017 by the author(s)

Tracing the path of forgetting in rule abstraction and exemplar retrieval

Janina A. Hoffmann

University of Konstanz

Bettina von Helversen

University of Basel

University of Zürich

Regina A. Weilbacher

University of Basel

Jörg Rieskamp

University of Basel

Word Count: 10'980

Abstract Count: 222

Author Note

This research was supported by Swiss National Science Foundation Grant 100014
_146169/1.

Correspondence concerning this article should be addressed to Janina A. Hoffmann,
Department of Psychology, University of Konstanz, Universitätsstrasse 10, 78 468
Konstanz, Germany. E-mail: janina.hoffmann@uni-konstanz.de

Abstract

People often forget acquired knowledge over time such as names of former classmates. Which knowledge people can access, however, may modify the judgment process and affect judgment accuracy. Specifically, we hypothesized that judgments based on retrieving past exemplars from long-term memory may be more vulnerable to forgetting than remembering rules that relate the cues to the criterion. Experiment 1 tracked the individual course of forgetting in a judgment task facilitating rule-based or exemplar-based strategies by systematically prolonging the retention interval between a training in which participants learned to make judgments and later tests (immediate, one day, and one week). Practicing the acquired judgment strategy in repeated tests helped participants to consistently apply the learnt judgment strategy and retain a high judgment accuracy even after a week. Yet, whereas a long retention interval did not affect judgments in the rule-based task, a long retention interval impaired judgments in the exemplar-based task. If practice was restricted as in Experiment 2, judgment accuracy suffered in both tasks. In addition, after a week without practice participants tried to reconstruct their judgments by applying rules in the exemplar-based task. These results emphasize that the extent to which decision makers can still retrieve previously learned knowledge limits their ability to make accurate judgments and that the preferred strategies change over time if the opportunity for practice is limited.

Keywords: Judgment, forgetting, rule-based and exemplar-based processes

Tracing the path of forgetting in rule abstraction and exemplar retrieval

One of the earliest discovered laws in psychology is the law of forgetting. The more time has passed between encoding an item and retrieving this item, that is the longer the retention interval is, the less likely people recall the item correctly (Ebbinghaus, 1885; Rubin & Wenzel, 1996). On a class reunion one year after high school, for instance, the names of former classmates may easily come to your mind. Twenty years later, however, you may even encounter problems when naming your former best friends. The course of time makes remembering facts, such as the names of previous classmates (Bahrick, Bahrick, & Wittlinger, 1975), or past events, such as headlines in newspapers (Meeter, Murre, & Janssen, 2005), more difficult.

If people forget information with the passage of time, this should also limit their ability to use this information when making judgments and decisions, affecting judgment quality. Although knowledge about how judgment accuracy varies as time passes by is limited (Ashton, 2000), it seems that not all judgments are equally affected by the time that has passed. For instance, meteorological forecasters have been shown to be more consistent than forecasters in the business or medical domain (Ashton, 2000). This domain difference could be due to people retrieving different information from memory depending on the judgment strategy they rely on.

Suppose, for instance, a judgment task where a juror in a song contest has to evaluate the performance of different singers on a scale from 0 to 10 every week. To judge the candidates, the jurors may consider how much the candidate meets different performance criteria, such as vocal skills or stage presence. If a juror has made up her mind how important the different performance criteria are and applies this policy consistently, her judgment should be independent of the time that passes between the candidates' performances or the number of performances she has heard before.

Alternatively, the juror may judge the candidate's performance by remembering how well past candidates performed on the show and how similar the current performance is to

these past performances. In this case the judgment will depend on how well the juror remembers these past performances. Accordingly, if the performances of previous candidates are remembered less well the more time has passed, judgments for a candidate should vary depending on the candidates that can still be retrieved.

In sum, understanding how judgment accuracy changes over time requires to understand which knowledge people retrieve when making a judgment and how knowledge retrieval changes with the passage of time. So far, however, this question has hardly been studied. The goal of the present research is to fill this gap and to investigate how the passage of time between learning a judgment task and making subsequent judgments influences judgment accuracy and interacts with the way people form their judgments. In the following we will describe the different judgment strategies people may follow and how they may be affected by forgetting in more detail.

Judgment strategies

Evaluating how well a singer performed in a song contest requires inferring a continuous criterion, the singers' performance, based on a number of features or cues of the judgment object, such as his vocal skills. People learn to solve these judgment tasks by getting feedback about the correct criterion. Judgment research has emphasized the idea that people can base such judgments on two types of judgment strategies: rule-based and exemplar-based (Erickson & Kruschke, 1998; Juslin, Karlsson, & Olsson, 2008, 2003; von Helversen & Rieskamp, 2008, 2009). These two strategies differ in the way they represent and process knowledge (Hahn & Chater, 1998; Juslin et al., 2003). Rule-based strategies assume that people try to test hypotheses about how each cue relates to the criterion (Brehmer, 1994; Juslin et al., 2008). To judge a new object, people integrate the weighted cue information linear additively (Einhorn, Kleinmuntz, & Kleinmuntz, 1979; Juslin et al., 2003). For instance, singers who more often strike the right chord and attract attention on stage may be evaluated more favorably. Accordingly, this judgment process

requires storing the importance assigned to each cue in long-term memory whereas information about previously encountered objects can be forgotten (Hoffmann, von Helversen, & Rieskamp, 2014; von Helversen & Rieskamp, 2009).

In comparison, exemplar-based judgment strategies assume that people store every previously encountered object, the exemplars, and their associated criterion values in long-term memory (Juslin et al., 2008, 2003; Nosofsky, 1988). To make a judgment, people retrieve all encountered objects from long-term memory and compare the current object (the probe) to all exemplars. The more similar the probe is to an exemplar, the more this exemplar influences the judgment. Hence, according to exemplar-based judgment strategies the juror stores each singer and their performance in long-term memory. The more a candidate's performance matches the best-rated candidates in previous shows, the more favorably the candidate will be evaluated.

Research suggests that people can adopt both kinds of judgment strategies, but shift between those strategies depending on the structure of the task (Juslin et al., 2008, 2003; Karlsson, Juslin, & Olsson, 2007; Pachur & Olsson, 2012; Platzer & Bröder, 2013; von Helversen & Rieskamp, 2009) and characteristics of the decision maker (Hoffmann et al., 2014; Little & McDaniel, 2015; von Helversen, Mata, & Olsson, 2010). Specifically, linear regression models, the predominant account to describe rule-based judgment strategies (Cooksey, 1996; Juslin et al., 2003), capture people's judgments well in linear judgment tasks in which the criterion is a linear, additive function of the cues. In contrast, exemplar models more accurately describe and predict participants' judgments in multiplicative tasks in which the criterion is a multiplicative function of the cues (Hoffmann et al., 2014, 2016; Juslin et al., 2008).

In sum, both rule-based and exemplar-based strategies require to some extent storage in and retrieval from episodic memory, but which kind of knowledge is stored and retrieved varies between the strategies. Whereas rule-based strategies assume that people need to store and retrieve each cue's importance, exemplar-based strategies rely on storage and

retrieval of past exemplars. Accordingly, both judgment strategies may be disrupted over time by forgetting, but forgetting may harm rules and exemplars to a different degree.

Sources of forgetting in rule-based and exemplar-based judgments

To what extent people forget information over time is a function of how well the information has been learned initially and if it can be successfully retrieved after some time has passed. The time that passed, however, may not cause forgetting per se, but rather what happened during this time (McGeoch, 1932; Rubin & Wenzel, 1996). Specifically, memory research has postulated two major mechanisms that may cause forgetting: a decay of stored memory traces and interference of similar items (for a historical review see Roediger, Weinstein, & Agarwal, 2010). Decay theories postulate that memory traces get weaker over time without accessing them. In contrast, interference theories postulate that storing similar items harms retrieval of the to-be-remembered items (Anderson & Neely, 1996; Postman, 1971). Accordingly, other memories compete with the target memory for retrieval and make it more difficult to retrieve the specific target item (Anderson, Bjork, & Bjork, 1994). Which mechanism underlies forgetting over a long time interval is hard to determine, but considering concepts from memory research may inform our understanding about how forgetting may affect retrieval of previously learned rules or exemplars.

Judgment tasks can be thought of as paired-associates learning tasks (Siegel & Siegel, 1972): During learning, people need to form an association between each cue and its importance in rule-based judgments, whereas they need to associate the exemplar (that is, a combination of cues) with its criterion value in exemplar-based judgments. During retrieval, the cues of the presented probe serve as retrieval cues for either the rule or the exemplar.

In rule-based judgments, it has been proposed that people abstract the importance of each cue and adjust (or update) its importance over trials (Hoffmann et al., 2014; Juslin et al., 2008; Pachur & Olsson, 2012). Once a participant has formed a satisfying rule, this

rule can be applied to each object. The established rule is hence practiced on every trial. Furthermore, the rule may generalize across different exemplars so that presenting a probe with a different combination of cues interferes with rule retrieval only to a small extent. As a result, rule-based judgments may not be harmed strongly by forgetting. Supporting this idea, Balzer, Rohrbaugh, and Murphy (1983) have found that judgments predicted from a rule-based regression model show a high test-retest reliability even after a week. Actual judgments, however, were less stable over time suggesting that forgetting still intrudes to some degree.

In exemplar-based judgments, it has been proposed that people store each exemplar in a separate memory trace (Estes, 1986). Frequently presented objects are more often encoded facilitating subsequent retrieval of the exemplar. Which exemplar is most similar to the probe, however, varies from trial to trial so that previously stored exemplars are practiced less often and may decay. Furthermore, stored exemplars likely share the same cue value on a particular cue so that the same cue value may activate exemplars with different criterion values. This overlap may increase competition between exemplars, disrupt discrimination between stored exemplars and, in turn, harm retrieval (Capaldi & Neath, 1995). In this vein, it has been found that people follow exemplar-based strategies less, if they cannot discriminate past exemplars from each other (Rouder & Ratcliff, 2004). In sum, exemplar-based judgment strategies may be more prone to forgetting than rule-based judgment strategies.

Previous studies indeed suggest that forgetting may harm retrieval of previously encountered exemplars more than retrieval of previously learned rules. In dot pattern classification paradigms, abstracted prototypes are better remembered over time than single instances (Homa, Cross, Cornell, Goldman, & Shwartz, 1973; Posner & Keele, 1970; Robbins et al., 1978). Furthermore, a recent study applying the looking-at-nothing paradigm found some evidence that people retrieve past exemplars more often in exemplar-based judgments than in rule-based judgments (Scholz, von Helversen, &

Rieskamp, 2015). In the looking-at-nothing paradigm (Richardson & Spivey, 2000) participants are presented with objects at different locations. During retrieval, participants tend to look back to the location at which the object they recall was presented suggesting that gaze location indicates which objects people retrieve. Using this looking-at-nothing paradigm, Scholz et al. (2015) found that people who base judgments on similarity look back more often to the location of previously seen similar exemplars than those who base judgments on rules.

If people do not remember previously encountered exemplars well after a long time interval, how can they still solve an exemplar-based judgment task? There is good reason to believe that people try to apply ill-conceived rules if previous exemplars cannot be retrieved (Olsson, Enkvist, & Juslin, 2006). In line with this idea, work relating memory abilities to judgment strategies has found that people with a better episodic memory more frequently adopt an exemplar-based strategy and, in turn, solve exemplar-based judgment tasks more accurately (Hoffmann et al., 2014). Furthermore, people who state that they relied on memory categorize new items more often based on similarity than those who indicated that they learned a rule (Little & McDaniel, 2015). Finally, Bourne, Healy, Kole, and Graham (2006) investigated how participants' stated classification strategy developed over the course of learning and changed after a one-week retention interval in different alphabetical categorization tasks. In the easy and difficult artificial tasks, participants indicated that rule use dominated early in learning, but over the course of learning more memory-based strategies evolved. After a week, however, participants stated that they relearned both tasks by applying a rule and, furthermore, did not revert to the memory-based strategy in the easy task. Accordingly, Bourne et al. (2006) argued that a longer retention interval induces a shift from memory-based to rule-based strategies because rules are better remembered than single instances.

Rationale of the current experiments

Taken together, both rule-based and exemplar-based judgments may involve to some extent storage in and retrieval from long-term memory: In rule-based judgment, people need to retrieve the previously learned rules. In exemplar-based judgment, they need to retrieve previously encountered exemplars and their criterion values. Rules are practiced on every trial and likely generalize across exemplars, whereas previously stored exemplars may be practiced less often and compete for retrieval. Accordingly, prolonging the retention interval between a training and a test phase may harm retrieval of single exemplars more than retrieval of rules.

One possibility to manipulate on which type of strategy people rely on is to vary the functional relationship between the cues and the criterion from a linear to a multiplicative one (Hoffmann et al., 2014, 2016; Juslin et al., 2008; Karlsson et al., 2007). In both task structures, participants judge the same objects, but the criterion value associated with each object varies between linear and multiplicative tasks. Linear tasks allow assessing the independent contribution of each cue to the criterion and thus testing linear rules is a viable strategy (Juslin et al., 2008). In comparison, multiplicative tasks require associating a combination of cues with the criterion value, but cannot be solved adequately by testing the independent effect of each cue so that participants should be more likely to rely on exemplar-based strategies (Juslin et al., 2008; von Helversen & Rieskamp, 2009). Consequently, we expected a stronger decline in judgment accuracy in multiplicative judgment tasks, which are more likely solved by exemplar memory than in linear judgment tasks in which people should predominantly try to abstract rules.

We tested this prediction in two experiments: Experiment 1 tracks the individual path of forgetting by asking participants to solve either a rule-based or an exemplar-based judgment task and repeatedly retrieve this knowledge: immediately, after a day, and after a week. In Experiment 2, we further explore the link between forgetting and judgment accuracy by manipulating the retention interval between participants from immediate recall

to recall after a week. Finally, we further tested to what degree forgetting may influence which cognitive strategies people tend to follow at each time point (Bourne et al., 2006).

Experiment 1: Forgetting over time with repeated practice

To test our hypotheses, we trained participants in the present study to predict the criterion value for a number of objects using four cues. In this training session, participants were randomly assigned to one of two judgment tasks: a linear judgment task to induce a rule-based judgment strategy or a multiplicative judgment task to induce an exemplar-based judgment strategy. To induce forgetting, we asked participants to judge old items (objects encountered in training) as well as new items (unknown objects) repeatedly at three retention intervals: an immediate test session, a test session after one day, and a test session after one week.

Method

Participants. 83 participants (53 female, 30 male, $M_{\text{Age}} = 24.6$, $SD_{\text{Age}} = 6.5$) were recruited at the University of Basel and randomly assigned to the linear ($n = 41$) or the multiplicative task ($n = 42$). Two participants who did not show up for all sessions were excluded from the study (one participant in the linear, and one in the multiplicative task) as well as one who was assigned to the wrong task in one of the sessions. Participants received course credit or 20 Swiss Francs (CHF) per hour for participating in the experiment. In addition, they could earn a bonus based on their judgment performance ($M = 6.06$ CHF, $SD = 2.11$ CHF). The first session took about an hour, whereas the second and the third lasted approximately 30 minutes.

Design and material. The cover story asked participants to predict how long the pupal stage lasts for different fictitious butterfly species on a scale from 10 to 20 days. The butterflies' appearance differed in four binary features (the cues): wing color (red vs. violet), antennae color (black vs. orange), color of stripes (brown or pink), and shape of spots (oval or serrated). Figure 1 shows two sample butterfly species with different cue

values on all cues. These cues could be used to predict how long the pupal stage for a butterfly lasts (the criterion). In the linear judgment task, the criterion was a linear, additive function of the cues,

$$y_{\text{lin}} = 4x_1 + 3x_2 + 2x_3 + x_4 + 10. \quad (1)$$

Each cue, x_1 , x_2 , x_3 , and x_4 , could take a cue value of zero or one. In the multiplicative judgment task, the criterion was a nonlinear, multiplicative function of the cues:

$$y_{\text{mult}} = 9 + e^{(4x_1+3x_2+2x_3+x_4)/4.15} \quad (2)$$

In principle, the exemplar model is able to learn to solve both types of judgment tasks, whereas a linear model only provides an accurate solution in linear tasks. Thus, to distinguish the models it is necessary to predict performance for new, unseen objects. From all possible items, we selected a training set of 10 old items and a test set of 6 new items so that a linear model and an exemplar model with one sensitivity parameter made different predictions for new items in both judgment tasks (see Appendix A for model descriptions). The sensitivity parameter determines whether people retrieve only highly similar exemplars or also more distant ones (Nosofsky & Zaki, 1998). In the linear judgment task, the linear model predicted the criterion values of new items more accurately than the exemplar model. In the multiplicative task, an exemplar model better fitted the old items and made slightly better predictions for new items than a linear model. Table 1 illustrates the task structure: The cues were given a binary value of zero or one, and they varied in their cue weights; that is, in their importance for predicting the criterion. The cue weights were randomly assigned to the four pictorial cues, as were the cue values (zero or one) to the features (e.g., oval or serrated spots).

One potential problem with manipulating the strategy with different task structures is that the task structure could also influence the ability to remember the training

exemplars. To ensure that the difference in model prediction did not depend on the task we simulated forgetting of exemplars in both tasks using an exemplar-based learning model (Kruschke, 1992). We modeled forgetting by assuming interference such that over time more and more new exemplars would be encountered that then would interfere with retrieving the previously learnt ones (for details see Appendix B). This exemplar-based learning model predicted that forgetting of exemplars would cause a similar increase in judgment error for old items in linear and multiplicative tasks, independent of the degree of interference. Repeating this simulation with parameters sampled from the best fitting parameters based on participants' training data led to the same conclusions, suggesting that the general predictions should not depend on the structure of the tasks we used.

Procedure. The judgment task consisted of a training session and three test sessions. In the training session, participants learned to estimate the criterion values for the 10 old items from the training set. In each trial, participants first saw a picture of a butterfly and were asked to estimate its criterion value. Afterwards they received feedback about the correct value, their own estimate, and the points they had earned. The training session ended when a learning criterion was reached. Participants met this learning criterion when judgment accuracy, as measured in root-mean-square deviation (RMSD) between participants' judgments and the criterion values in one training block, fell below 1 RMSD. We employed this learning criterion to minimize the possibility that differential forgetting in judgment could have resulted from initial differences in judgment accuracy between tasks and to achieve a high judgment accuracy at the end of training. Each participant completed at least 10 training blocks, each consisting of the 10 old items; training terminated after 20 blocks even if the learning criterion had not been reached. Earlier work has set a more lenient learning criterion of 1.5 RMSD to be met within fewer training blocks (Mata, von Helversen, Karlsson, & Cüpper, 2012; von Helversen et al., 2010) suggesting that participants may meet our learning criterion as well within 20 training blocks. In the test sessions, participants estimated criterion values for all 16

butterflies, 10 old and 6 new ones, six times without getting any feedback. Presentation order in each training and test block was randomly determined.

To motivate participants to reach a high performance, participants could earn points in every trial. Participants earned 10 points for a correct answer and 5 points if their judgment deviated by 1 from the correct answer. At the end of the judgment tasks, the points earned were converted to a monetary bonus (500 points = 1 CHF). In addition, participants earned a bonus of 5 CHF if they reached the learning criterion for the judgment task within 20 training blocks. Participants returned to the lab after 24 h as well as after one week to repeat the test session of the judgment task.

Results

Learning success at the end of training. Overall, the number of participants reaching the learning criterion varied slightly between the judgment tasks, but the difference was not significant, $\chi^2(1) = 2.20, p = .138$. In the linear judgment task, 25 out of 40 participants reached the learning criterion (62.5%), whereas in the multiplicative task 32 out of 40 participants (80%) mastered the training phase successfully. Among those participants who did not learn the task, three participants in the multiplicative task and four participants in the linear task did not outperform a guessing model (a model predicting the mean of the criterion value in every trial). In the linear judgment task, however, participants needed slightly more training blocks ($M = 15.3, SE = 0.7$) than in the multiplicative task, $M = 13.8, SE = 0.7, t(78) = 1.6, p = .108$. The number of training blocks participants needed was highly correlated with judgment error in both tasks (linear: $r = .75$; multiplicative: $r = .78$).

Judgment performance over time. According to our hypothesis, increasing the retention interval between training and test should increase judgment error more in the multiplicative than in the linear judgment task. This increase in judgment error should be most pronounced for old items because people should be more likely to forget specific

training exemplars than previously learned rules. To compare judgment performance for old items across time, we measured judgment error in the training session as the RMSD between participants' judgments in the last training block and the correct criterion. Judgment error in the three test sessions (immediate test, test after 1 day, and after 1 week) was measured as the RMSD between the criterion and participant's judgments, averaged over the six presentations in each test session. Figure 2 shows judgment error for old and new items in each test session separately for the linear (white bars) and the multiplicative judgment task (gray bars). Descriptively, participants achieved a similar accuracy level for old items at the end of training in the linear (RMSD = 1.23, $SE = 0.17$) and the multiplicative judgment task (RMSD = 1.02, $SE = 0.16$). Judgment error in the linear task is equally high (even slightly lower) as in studies using a similar design (Mata et al., 2012; von Helversen et al., 2010). In the linear judgment task, judgment error remained constant across time from immediate test (RMSD = 1.16, $SE = 0.15$) to the next day (RMSD = 1.20, $SE = 0.14$) to one week (RMSD = 1.19, $SE = 0.14$, $d = -0.08$ from last block of training to one week using an effect size based on the change score for repeated measures, Morris & DeShon, 2002). In the multiplicative task, however, judgment error rose for old items from immediate test (RMSD = 1.28, $SE = 0.15$) to the next day (RMSD = 1.45, $SE = 0.16$) to a week later (RMSD = 1.55, $SE = 0.16$, $d = 0.70$ from last block of training to one week).

To test the hypothesis that a longer retention interval harms exemplar-based judgments more than rule-based judgments on old training items we conducted a linear mixed model analysis using judgment error for old items as the dependent variable and compared models with different predictors using a likelihood ratio test. Further, we report Akaike's Information criterion (AIC) for each model (Janssen, 2012). Session (training, immediate test, test after 1 day, and test after 1 week) and judgment task (linear vs. multiplicative) were included as fixed factors, participants and items as random intercepts (random model: AIC = 9738). The type of judgment task did not affect judgment error,

AIC = 9740, $\chi^2(1) < 0.1$, $p = .953$, but participants made less accurate judgments in later test sessions, AIC = 9720, $\chi^2(3) = 26.4$, $p < .001$. Furthermore, session interacted with the type of judgment task, AIC = 9709, $\chi^2(3) = 16.6$, $p < .001$, suggesting that judgment error increased more over time in the multiplicative than in the linear judgment task. To further investigate in which sessions judgment error increased most strongly in the two tasks, we analyzed both tasks separately using contrasts on the least square means for the different sessions. In the multiplicative task, comparing judgment error at the end of training to average judgment error in the three test sessions suggested that judgments at the end of training were more accurate than judgments in test, $b = 1.11$, $SE = 0.18$, $t(3105) = 6.06$, $p < .001$. Furthermore, participants made fewer errors in immediate test than in the delayed tests, $b = 0.30$, $SE = 0.13$, $t(3105) = 2.32$, $p = .020$, but judgment error did not increase from delayed test after a day to test after one week, $b = 0.05$, $SE = 0.07$, $t(3105) = 0.68$, $p = .497$. In the linear judgment task, judgment error increased neither from the last block of training to the three test sessions, $b = 0.14$, $SE = 0.18$, $t(3105) = 0.76$, $p = .446$, nor from immediate test to delayed tests, $b = 0.02$, $SE = 0.13$, $t(3105) = 0.16$, $p = .872$, nor from one day to one week, $b = -0.00$, $SE = 0.07$, $t(3105) = -0.06$, $p = .956$.¹

With regard to new items, judgment error descriptively increased in the multiplicative judgment task from immediate test (RMSD = 2.50, $SE = 0.17$) to delayed test after a day (RMSD = 2.54, $SE = 0.17$) and test after a week (RMSD = 2.68, $SE = 0.18$, $d = 0.26$ from immediate to one week). Likewise, judgment error increased for new items in the linear judgment task from immediate test (RMSD = 1.56, $SE = 0.15$) to test

¹Excluding participants based on the learning criterion slightly changed results suggesting in addition that participants in the linear task performed better across all sessions, $\chi^2(1) = 4.2$, $p = .040$. Inspecting the interaction more closely further indicated that, judgment error increased in the linear task as well from training to test, $b = 0.50$, $SE = 0.18$, $t(2208) = 2.82$, $p = .005$, but neither between immediate and delayed tests, $b = 0.13$, $SE = 0.13$, $t(2208) = 1.07$, $p = .285$, nor between test after a day and after a week, $b = -0.01$, $SE = 0.07$, $t(2208) = -0.07$, $p = .942$. The pattern for strategy use remained the same so that so that we report results for all participants.

after one day ($RMSD = 1.69$, $SE = 0.18$) to test after one week ($RMSD = 1.77$, $SE = 0.19$, $d = 0.33$ from immediate to one week). However, a mixed model analysis with participants and items as random intercepts (random model: $AIC = 5011$) indicated that participants made worse judgments on new items in the multiplicative judgment task than in the linear judgment task, $AIC = 5002$, $\chi^2(1) = 10.7$, $p = .001$, but did not make significantly more errors in later test sessions, $AIC = 5001$, $\chi^2(2) = 4.6$, $p = .102$. Furthermore, there was no interaction between session and type of judgment task, $AIC = 5005$, $\chi^2(2) = 0.6$, $p = .753$.

Taken together, we found that prolonging the retention interval between training and test increased judgment errors on old training items in the multiplicative judgment task, but not in the linear judgment task. If participants had to generalize their judgment to new items, we found that participants in the multiplicative judgment task made on average more errors than participants in the linear judgment task. A longer retention interval, however, did not affect judgment errors on new items.

Judgment strategies over time. Previous research (Bourne et al., 2006) has suggested that a long retention interval leads to a shift from exemplar-based strategies to rule-based strategies because people cannot retrieve previously encountered exemplars and instead they relearn the task by applying a rule. According to this hypothesis, participants in the multiplicative judgment task should shift from exemplar-based strategies in the immediate test to rule-based strategies after a week. However, in contrast to Bourne et al. (2006), in our study participants had the possibility to repeatedly practice their judgment strategy both in the immediate test and after a day, making it likely that they did not need to abandon an exemplar-based judgment strategy after a week. To classify participants as following a rule-based or exemplar-based strategy, we relied on a cognitive modeling approach. We fitted a linear regression model serving as a rule-based model, and an exemplar model with four attention weights to participants' judgments, separately for each of the three test sessions (see Appendix A for more details on modeling and Table A1 for participants' strategy classifications). To account for random guessing, we compared those

models to a baseline model (a model estimating participants' mean judgment). Three participants who were best described by a baseline model in one or more sessions were excluded from all following analyses (two participants in the linear and one in the multiplicative task). In the linear judgment task, most participants were best described by a rule-based model in all three test sessions (immediate: 32, one day: 29, one week: 34), whereas only a few participants were best described by an exemplar model (immediate: 6, one day: 9, one week: 5). In the multiplicative task, the number of people classified to the exemplar model increased compared to the linear task (immediate: 15, one day: 16, one week: 14). However, a high number of participants still followed a rule-based model (immediate: 24, one day: 23, one week: 26).

To assess how closely the judgment strategies in immediate test corresponded to the judgment strategies after a day or after a week, we calculated Cohen's κ using the percentage of concordant classifications and corrected for the number of categories because the number of observations per category was not equal (Brennan & Prediger, 1981). In the multiplicative task, the judgment strategy that participants relied on in immediate test corresponded closely to the strategy participants relied on after a day, Cohen's $\kappa = 0.74$, and to the strategy participants relied on after a week, Cohen's $\kappa = 0.70$. In the linear task, agreement was similarly high between judgment strategies in immediate test and a day later, Cohen's $\kappa = 0.74$, and slightly lower a week later, Cohen's $\kappa = 0.66$.

To investigate to what extent participants shifted between strategies as a function of the retention interval, we conducted a logistic mixed-model regression using judgment strategy as the dependent variable and test session as well as judgment task as independent fixed factors with participant as a random intercept (AIC = 227). A main effect of judgment task, AIC = 221, $\chi^2(1) = 7.30$, $p = .007$, indicated that more participants were better described by exemplar models in the multiplicative than in the linear task, odds ratio (OR) = 13.9, 95% confidence interval (CI) [2.1, 172.9]. However, there was neither a main effect of session, AIC = 222, $\chi^2(2) = 3.71$, $p = .156$, nor an interaction between

session and judgment task, $AIC = 224$, $\chi^2(2) = 1.37$, $p = .504$. In sum, the type of judgment task predicted which judgment strategy described participants' judgments best at each time point, but a longer retention interval did not increase the number of participants best described by rules suggesting that participants did not shift to rule-based judgment strategies in response to a longer retention interval.

Discussion

In Experiment 1, we investigated whether prolonging the retention interval affects how accurately people make judgments in two different kinds of judgment tasks: a linear judgment task that can best be solved by abstracting linear, additive rules and a multiplicative judgment task that can better be approached by storing and retrieving exemplars from long-term memory. In line with our hypothesis, we found that judgment accuracy for old items encountered in training dropped more from training to recall after a week in the multiplicative than in the linear judgment task. In the linear judgment task, participants judged —on average —old items as accurately after a week as at the end of training, whereas judgment errors increased from the last training block to test after a week in the multiplicative judgment task. This result matches previous research suggesting that people remember abstracted knowledge, for instance in the form of prototypes, better than single instances after a long retention interval (Homa et al., 1973; Posner & Keele, 1970; Robbins et al., 1978) and supports the idea that exemplar-based judgments build to a stronger extent on episodic memory than rule-based judgments (Hoffmann et al., 2014).

Replicating previous work on judgment strategies (Hoffmann et al., 2014; Juslin et al., 2008), we found that more participants were best described by an exemplar model in the multiplicative task than in the linear task. With regard to the question of how judgment strategies developed over time, our results suggest that participants relied consistently on the same judgment strategy across time: In the linear task, most participants were still best described by a rule-based strategy after a week. Similarly, the same number of

participants in the multiplicative task was best described by an exemplar model after one week suggesting that those participants still tried to retrieve previously encountered exemplars. This finding differs from previous research (Bourne et al., 2006) suggesting that people prefer relearning complex categorizations by relying on rules, although they stated that they previously solved the task by retrieving exemplars from memory.

One reason why participants potentially did not shift from an exemplar-based strategy to a rule-based strategy is that they had the possibility to repeatedly practice their judgment strategy. Practicing a task even without getting feedback can benefit long-term retention in a wide range of tasks from free recall to function learning and may outperform studying the correct solution (Kang, McDaniel, & Pashler, 2011; Karpicke & Roediger, 2008). One explanation why practice is so beneficial for retention focuses on the idea that those repeated retrieval processes may strengthen the memory trace by elaboration, deeper encoding or adding multiple retrieval routes (Roediger & Karpicke, 2006). It is possible that asking participants to solve the judgment task immediately, after one day, *and* after one week, involved such repeated retrieval processes. Therefore, our design that tried to track individual paths of forgetting might have prevented a high amount of forgetting in the judgment task. To circumvent the possibility that repeated practice may have restricted the amount of forgetting, we tested in a second experiment whether forgetting impacts judgments more and participants shift to a greater extent to rule-based strategies if they do not have the opportunity to repeatedly practice the judgment task at several time points.

Experiment 2: Forgetting over time without repeated practice

In Experiment 2, we studied how forgetting affected participants' judgments if they did not have the possibility to practice their judgment strategy between training and a later test. As in Experiment 1, participants learned to solve either a linear or a multiplicative judgment task in a training session. To induce forgetting, we asked participants to rejudge these old items as well as new items either immediately after training or after one week.

In addition, we assessed recognition memory for old items in a two-alternative forced-choice test. Past research has found that participants who possess a better episodic memory are more likely to adopt an exemplar-based strategy and, in turn, make more accurate judgments in multiplicative judgment tasks (Hoffmann et al., 2014). This suggests that people using an exemplar strategy may remember better which objects they encountered during training than rule users. However, there is some research suggesting that the relationship between recognition and strategy use is more complex. First, if people are asked to recognize the previously encountered exemplars in a recognition test they are better at remembering items violating a salient knowledge structure, for instance a rule, than items following the knowledge structure (Davis, Love, & Preston, 2012; Palmeri & Nosofsky, 1995; Sakamoto & Love, 2004). Second, the more salient the rules are in these rule-plus-exception tasks, the better the exceptions violating the rule are remembered (Sakamoto & Love, 2004). This result could indicate that also rule users can perform well in a multiplicative task if they can remember the exceptions to the rule well. Lastly, it has been found in rule-plus-exception tasks that previously encountered items consistent with a rule show a recognition advantage over novel items implying that people possess some residual memory for old exemplars, although they abstracted a rule (Palmeri & Nosofsky, 1995). Accordingly, it is possible that learners in the rule-based judgment task also encode an exemplar-specific representation and show a recognition advantage for previously encountered old items over new items reducing differences in recognition performance between strategies and tasks.

Method

Participants. 142 participants (115 female, 27 male, $M_{\text{Age}} = 24.3$, $SD_{\text{Age}} = 6.4$) were recruited at the University of Basel. Participants were randomly assigned to one of the four conditions: 35 to the linear task with immediate recall, 37 to the linear task with recall after a week, 33 to the multiplicative task with immediate recall and 37 to the

multiplicative task with recall after a week. Two participants who did not show up for the test session after one week were excluded from the study (one participant in the linear, and one in the multiplicative task) as were three participants who were assigned to the wrong condition. Participants received course credit or 20 Swiss Francs (CHF) per hour for participating in the experiment. In addition, they could earn a performance-dependent bonus ($M = 5.17$ CHF, $SD = 2.68$ CHF). The training session took about an hour, whereas the test session took approximately thirty minutes.

Procedure. Material and procedure followed closely the experimental set-up of Experiment 1. Participants were randomly assigned to the linear or the multiplicative judgment task. In contrast to Experiment 1, we varied the retention interval between training and test between participants: Half of the participants in each judgment task solved the test session immediately after training whereas the other half returned to the lab after a week.

After participants completed the test session, they solved a two-alternative forced-choice recognition test. In each trial, participants saw one "old" butterfly—that is, one they already knew from training—and one "new" butterfly—that is, a butterfly from the test set introduced in the test session. Participants had to determine which of those two butterflies was "old"; that is, the one they already knew from training. All 10 old butterflies were presented twice with each of the 6 new butterflies, resulting in 120 forced choice decisions.

Results

Learning success at the end of training. As in Experiment 1, the number of participants reaching the learning criterion did not vary strongly between the judgment tasks and the retention intervals, $\chi^2(4) = 3.7$, $p = .448$. In the linear judgment task, 20 out of 35 participants (57.1%) assigned to immediate test reached the learning criterion as did 19 out of 36 participants (52.8%) assigned to test after a week. The multiplicative task was

mastered successfully by 23 out of 31 participants (74.2%) assigned to immediate test and by 21 out of 35 participants (60.0%) assigned to the test one week later. Among those participants who did not learn the task, five participants in the linear task (immediate: 2, one week later: 3) and 13 participants in the multiplicative task (immediate: 4, one week later: 9) did not outperform a random guessing model. Table 2 displays descriptive statistics for both judgment tasks, separately for immediate test and test after one week. In the multiplicative task, participants needed —on average— slightly fewer training blocks than participants in the linear judgment task. In a two-way analysis of variance on the number of training blocks, neither the type of judgment task affected the number of blocks, $F(1,133) = 3.31$, $p = .071$, nor did the retention interval, $F(1,133) = 2.32$, $p = .130$, nor their interaction, $F(1,133) = 0.01$, $p = .942$. As in Experiment 1, judgment error in the last training block and the number of training blocks needed were highly correlated ranging from $r = .69$ in the linear task for test after a week to $r = .77$ in the linear task for immediate test.

Judgment performance over time. As in Experiment 1, we expected a longer retention interval to impede judgment accuracy most severely on old items in the multiplicative judgment task. Figure 3 illustrates judgment error on old and new items for the last block of training and the test session, separately for the judgment tasks and retention intervals (see Table 2 for descriptive statistics). In the linear judgment task, participants who took the immediate test were descriptively as accurate on old items in test as in the last block of training ($d = -0.19$, d based on the change score for repeated measures Morris & DeShon, 2002), whereas participants who solved the test session a week later made more errors on old items in test than in the last block of training ($d = 0.21$). In the multiplicative judgment task, participants who took the immediate test made only slightly more errors on old items in test than in the last block of training ($d = 0.23$), whereas participants who solved the test session a week later made worse judgments on old items in test than in the last block of training ($d = 0.64$). Finally, participants who solved

the linear judgment task a week later made on average fewer errors than participants who solved the multiplicative task, $d = -0.61$.

To test the hypothesis that a longer retention interval harms exemplar-based judgments more than rule-based judgments for old items we conducted a linear mixed model analysis on judgment errors using retention interval, judgment task, and session (training vs. test) as fixed factors and participant as well as item as random intercepts (random model: AIC = 9768). Overall, participants made fewer errors in the last training block than in test, AIC = 9728, $\chi^2(1) = 42.7$, $p < .001$, but the judgment task did not affect judgment errors, AIC = 9728, $\chi^2(1) = 1.2$, $p = .263$. A longer retention interval increased judgment error, AIC = 9722, $\chi^2(1) = 8.0$, $p = .005$. An interaction between retention interval and session indicated that judgment error increased more strongly between training and test for those participants who took the test after a week than immediately, AIC = 9707, $\chi^2(1) = 17.8$, $p < .001$. Further, an interaction between judgment task and session indicated that judgment error increases more from training to test for participants in the multiplicative task than for participants in the linear task, AIC = 9698, $\chi^2(1) = 10.6$, $p = .001$. Yet, in contrast to our hypothesis that a longer retention interval contributes to more errors in the multiplicative than in the linear judgment task, neither the interaction between retention interval and judgment task, AIC = 9700, $\chi^2(1) = 0.6$, $p = .458$, nor the three-way interaction was significant, AIC = 9701, $\chi^2(1) = 1.0$, $p = .319$.

As in Experiment 1, we broke up the interactions by separately analyzing the judgment tasks. In the multiplicative task, judgment errors increased both from training to test for participants who took a test after a week, $b = 0.77$ for the mean difference, $SE = 0.10$, $t(2590) = 7.6$, $p < .001$, as well as for participants who took the immediate test, $b = 0.23$, $SE = 0.11$, $t(2590) = 2.2$, $p = .030$. Comparing the size of the increase suggests a stronger increase in error for participants who took the delayed than the immediate test, $b = 0.54$, $SE = 0.15$, $t(2590) = 3.6$, $p < .001$.

In the linear task, participants in immediate test did not make more errors in test than in training, $b = 0.06$, $SE = 0.10$, $t(2590) = 0.06$, $p = .955$, whereas participants who took a test after a week made worse judgments in test than at the end of training, $b = 0.34$, $SE = 0.10$, $t(2590) = 3.3$, $p < .001$. As in the multiplicative task, the increase in error was more pronounced for participants who took the delayed than the immediate test, $b = 0.33$, $SE = 0.14$, $t(2590) = 2.3$, $p = .021$.²

With regard to new items, participants in the linear task descriptively made more errors if they were tested a week later than if they took an immediate test ($d = 0.13$). Similarly, in the multiplicative task participants who were tested after a week made less accurate judgments than those who took an immediate test ($d = 0.41$). Furthermore, participants made less accurate judgments in the multiplicative task than in the linear task both in immediate test ($d = 0.57$) and after a week ($d = 0.79$). To investigate how a longer retention interval affected judgment errors for new items we conducted a linear mixed model analysis on judgment errors using retention interval and judgment task as fixed factors and item as random intercept (random model: AIC = 3090). Judgments were less accurate in the multiplicative compared to the linear task, AIC = 3067, $\chi^2(1) = 25.4$, $p < .001$, but neither retention interval, AIC = 3066, $\chi^2(1) = 3.1$, $p = .079$, nor its interaction with the type of task affected judgment accuracy, AIC = 3065, $\chi^2(1) = 3.0$, $p = .083$.

In sum, a longer retention interval impeded judgment accuracy on old items in both judgment tasks. Judgment error increased more from training to test for those participants who took a delayed test after a week than for those who took an immediate test. Furthermore, as in Experiment 1, participants in the multiplicative judgment task were less successful at generalizing their performance to new items than participants in the linear task, independent of the retention interval.

²Excluding participants based on the learning criterion only slightly changed results for judgment error on training items. Specifically, this analysis suggested a significant three-way interaction between judgment task, session, and retention interval, $\chi^2(1) = 5.0$, $p = .025$, but the pattern of results remained the same within each judgment task. The analysis on strategy use did not yield a different pattern of results.

Judgment strategies over time. In Experiment 2, participants did not have the possibility to practice their judgment strategy between training and delayed test. Without practicing the judgment strategy, it is possible that participants shift from an exemplar-based judgment strategy to a rule-based judgment strategy after a week (Bourne et al., 2006). To describe judgment strategies, we fitted an exemplar model, a rule-based model and a baseline model to participants' judgments in each test session. In all subsequent analyses, we excluded participants best described by the baseline model: two participants in the linear task (immediate: $n = 1$, one week: $n = 1$) and six participants in the multiplicative task (immediate: $n = 2$, one week: $n = 4$). In the linear task, most participants were best described by a rule-based model in immediate test ($n = 27$, exemplar: $n = 7$) as well as in a test after a week ($n = 29$, exemplar: $n = 6$). In the multiplicative task, the exemplar model described more participants best in immediate test ($n = 18$) than the rule-based model ($n = 11$). After a week, however, only 8 participants were best described by an exemplar model whereas the majority of participants was best described by a rule-based model ($n = 23$).

We conducted a logistic regression analysis to understand how strategy use changes depending on the type of judgment task and retention interval. A model using retention interval and judgment task as predictors predicted strategy use best (Deviance $D = 143$, Nagelkerke's $R^2 = .156$, Cox and Snell's $R^2 = .110$) and outperformed a model relying only on judgment task, $D = 149$, Nagelkerke's $R^2 = .098$, Cox and Snell's $R^2 = .069$, $\chi^2(1) = 5.8$, $p = .016$. Adding an interaction between retention interval and judgment task did not improve model fit, $D = 141$, Nagelkerke's $R^2 = .180$, Cox and Snell's $R^2 = .127$, $\chi^2(1) = 2.5$, $p = .113$. The best fitting model suggests that fewer participants relied on an exemplar model after a week than in immediate test, $OR = 0.38$, $CI = [0.17; 0.68]$. Furthermore, in the multiplicative task more participants were best described by the exemplar model than in the linear task, $OR = 3.5$, $CI = [1.6; 8.2]$. Taken together, those results indicate that participants may have shifted from a memory-based exemplar strategy in immediate test to

a rule-based strategy after a week.

Predicting recognition memory with judgment strategies. Finally, we assessed in both tasks to what extent strategy use can predict recognition memory for previously encountered exemplars across time. On the one hand, participants relying on an exemplar-based strategy may rely more on episodic memory and discriminate old from new items better than participants relying on rules. On the other hand, it is possible that also rule-based learners may possess some residual memory for old exemplars (Palmeri & Nosofsky, 1995) and are likewise able to discriminate old from new items. Overall, participants correctly recognized 63.0% (*recognition rate*, $SD = 18\%$) of all old items. In the linear judgment task, the recognition rate was higher in test after a week than in immediate test (see Table 2). In the multiplicative task, participants recognized the old items slightly worse after a week than in immediate test.

To investigate how judgment strategies and retention interval affected recognition memory, we conducted a mixed logistic regression using the number of correctly and incorrectly recognized old items as dependent variable. Retention interval, judgment task, and judgment strategy were included as fixed factors in the analysis, whereas item was included as a random intercept (random model: $AIC = 8119$). Overall, the recognition rate was higher in the linear than in the multiplicative judgment task, $AIC = 8081$, $\chi^2(1) = 40.7$, $p < .001$. Furthermore, this analysis suggested a three-way interaction between judgment task, retention interval, and judgment strategy, $AIC = 7951$, $\chi^2(1) = 31.7$, $p < .001$. Therefore, we broke up this interaction by separately analyzing the judgment tasks. In the linear task (random model: $AIC = 4333$), participants had a higher recognition rate after a week than in immediate test, $AIC = 4286$, $\chi^2(1) = 49.2$, $p < .001$. Furthermore, participants classified to the rule-based model recognized more old items correctly than participants classified to the exemplar model, $AIC = 4249$, $\chi^2(1) = 38.5$, $p < .001$, and retention interval interacted with the judgment strategy, $AIC = 4241$, $\chi^2(1) = 10.5$, $p = .001$. This interaction suggested that participants classified to the rule-based model already

recognized old items better in immediate test than exemplar users, $b = .19$, $SE = .08$, but recognition performance further increased compared to exemplar users in the delayed test, $b = .57$, $SE = .09$. In the multiplicative task (random model: $AIC = 3760$), participants overall had a lower recognition rate after a week than in immediate test, $AIC = 3755$, $\chi^2(1) = 6.1$, $p = .013$, and the recognition rate differed between participants classified to the rule-based and the exemplar model, $AIC = 3743$, $\chi^2(1) = 14.8$, $p < .001$. Furthermore, an interaction between retention interval and judgment strategy, $AIC = 3722$, $\chi^2(1) = 23.1$, $p < .001$, indicated that recognition did not differ between participants classified to the rule and exemplar model in immediate test, $b = -.03$, $SE = .07$, but participants classified to the exemplar model had a higher recognition rate after a week, $b = .48$, $SE = .08$.

Figure 4 shows the recognition rate for each old item, plotted over participants' average judgment for this item. The graph illustrates that participants who are classified to the rule-based model in the linear task recognize old items better than participants who are classified to the exemplar model. In the multiplicative task, after a week participants classified to the exemplar model recognize old items better than participants who are classified to the less suitable strategy. However, exemplar users also show larger standard errors than rule-users after a week indicating that recognition memory has a higher variability. The reason for this finding is possibly that only a few participants still adopt an exemplar-based strategy after one week.

Discussion

Instead of tracking the individual course of forgetting, Experiment 2 varied the retention interval between participants to reduce the possibility that repeated practice of judgment strategies limited the decline of judgment accuracy over time. In line with the results from Experiment 1, we found that participants in the multiplicative task judged old items less accurately in immediate test than in training, whereas participants in the linear judgment task kept their performance in immediate test. In contrast to our hypothesis,

however, judgment error increased more strongly in a delayed test after a week not only in the multiplicative task, but also in the linear judgment task. Accordingly, a longer retention interval harmed judgments both in the multiplicative and in the linear judgment task. Furthermore, in contrast to Experiment 1, judgment strategies were not stable over time, but changed across time: In immediate test, the majority of participants in the linear task was classified to a rule-based model, whereas more participants were classified to an exemplar model in the multiplicative task. After a week, however, participants relied less on an exemplar model in both judgment tasks suggesting that if participants do not have the opportunity to practice an exemplar-based judgment process, they shift to a greater extent to rule-based strategies.

General discussion

The passage of time makes it harder to remember previously learned knowledge (Ebbinghaus, 1885; Rubin & Wenzel, 1996), but it can also impede previously acquired skills, such as speaking foreign languages (Bahrick, 1984). Although forgetting affects a wide range of cognitive abilities, only a few studies in judgment research have paid attention to such basic memory phenomena. Our research tried to shed light on the question of how a longer retention interval may change the knowledge people retrieve to make a judgment and, ultimately, judgment accuracy. Reinterpreting judgment tasks as paired-associates learning tasks, we argued that people may need to form different associations when learning to solve rule-based and exemplar-based judgment tasks: In rule-based judgments, people should associate each cue with its importance, whereas they need to associate exemplars with their corresponding criterion value in exemplar-based judgments. In a later test phase, people retrieve either previously learned rules or exemplars. Specifically, we hypothesized that storing a range of similar exemplars may make exemplar-based judgments highly vulnerable to forgetting, whereas rules receive more training, are likely generalized to a range of different objects, and may hence be forgotten

less easily. We tested this hypothesis in two experiments: one that tracked participant's judgment performance after different retention intervals and one that varied retention interval between groups. In both experiments, we found that judgment error on old items increased more from training to test in the multiplicative than the linear judgment task, reflecting the idea that forgetting over time harms successful retrieval of single exemplars more than retrieval of rules.

Looking more closely at the temporal curve of forgetting, we found that judgment error for previously encountered items in the multiplicative task already increased between the end of training and immediate test in both experiments, whereas participants in the linear judgment task kept their performance from training to immediate test. Possibly, introducing new items in the multiplicative task already interferes with retrieving old training exemplars so that participants likely confuse old training items with novel items. Yet, our results provide mixed evidence for the idea that prolonging the retention interval leads to a greater amount of forgetting in exemplar-based than in rule-based judgment. In line with previous research suggesting that forgetting does not act on abstracted knowledge like prototypes (Homa et al., 1973; Posner & Keele, 1970; Robbins et al., 1978) we found in Experiment 1 that participants in the linear judgment task were able to retain a high judgment accuracy even after a week, whereas participants in the multiplicative task made more errors in the delayed tests. In Experiment 2, however, a delayed test harmed judgment accuracy to the same degree in the multiplicative, exemplar-based judgment task as in the linear, rule-based judgment task. This result matches previous findings suggesting that actual rule-based judgments can also fluctuate over time (Balzer et al., 1983), but stands in contrast to research suggesting that abstracted knowledge is immune to forgetting (Homa et al., 1973; Posner & Keele, 1970; Robbins et al., 1978). One reason why people better retain rule-based judgments in Experiment 1 is possibly that they were able to apply the rules learned in training immediately to new test items so that the learned rules are less vulnerable to forgetting in the delayed tests. Limiting this opportunity to

practice the judgment strategy, as in Experiment 2, may have restricted not only repeated retrieval of exemplars, but also the generalization of rules to new items. Taken together, those results point towards the view that not only exemplars may be forgotten over a longer time interval, but people may also experience difficulties to retrieve previously learned rules after a long time. Future research may seek to unravel on a more fine-grained level the degree to which specific mechanisms of forgetting, such as decay or interference, underlie forgetting in rule and exemplar retrieval.

When participants had to generalize the learned knowledge to new items, we found that participants did not make more errors over time in both judgment tasks, but participants in the multiplicative task were rather bad at judging new items independent of the retention interval. A plausible reason for this high number of judgment errors is that we selected the new items so that they strongly discriminate between the strategies, but both strategies did not generate a high performance on new items in the multiplicative task. Accordingly, we used the new items primarily to distinguish the judgment strategies and not to evaluate performance.

The stability of judgment strategies and exemplar memory over time

The question of how stable people's judgment strategies are over time is of high practical relevance (Ashton, 2000). One line of research has argued that people's judgment weights may fluctuate only to a small degree over time (Balzer et al., 1983), whereas other researchers have proposed that the time that has passed critically influences the strategy people follow (Bourne et al., 2006). Our study unites those divergent ideas: If participants had the opportunity to repeatedly practice their judgment strategy, we found that their judgment policies were highly consistent across time indicating that repeated practice can render people's judgment policies temporally more stable. However, if participants did not engage in exemplar retrieval for a long time as in Experiment 2, they shifted more towards rule-based strategies. These findings are consonant with the idea that people only engage

in exemplar retrieval after a long time if the exemplars can still be retrieved. However, if previously encountered exemplars can no longer be retrieved, participants may revert to a less appropriate rule-based strategy (Bourne et al., 2006; Olsson et al., 2006).

To assess to what extent people still possess some memory for specific exemplars after a week, we additionally measured recognition memory in Experiment 2. Interestingly, the ability to discriminate old from new items varied in both judgment tasks as a function of the retention interval and strategy used: In the multiplicative task, exemplar and rule users were equally successful in discriminating between old and new items in immediate test; however, a week later, participants classified to the exemplar model more accurately discriminated between old and new items than participants classified to the rule-based model—a finding further supporting the idea that those participants who have a worse memory for previously encountered exemplars try to reinstate their judgment by applying rules. In turn, in the linear task, participants who were best described by the task-appropriate rule-based strategy better recognized old items than participants best described by the exemplar model and this recognition advantage increased in delayed test. Furthermore, on average, participants more accurately discriminated between old and new items in the linear than the multiplicative task. In combination, these findings highlight that rule-based learners also possess some residual memory for each exemplar (Palmeri & Nosofsky, 1995; Sakamoto & Love, 2004).

Restrictions in training duration

In our study, we tried to equate learning performance by setting a strict learning criterion, but participants could achieve this learning criterion after a variable number of learning blocks. We used this learning criterion because participants solving a multiplicative task often achieve a higher performance than participants in the linear task after the same number of training blocks (Hoffmann et al., 2014, 2016). In our study, participants in the multiplicative task also reached the learning criterion slightly faster

than participants in the linear task. This result may hint at the alternative interpretation that training may have prematurely stopped in the multiplicative task and therefore participants may have remembered their judgments less well. Two arguments speak against this interpretation. First, if participants in the multiplicative task reached the learning criterion by chance, they should make more errors in the blocks preceding the last training block than participants who passed the learning criterion in the linear task. Yet, in the three blocks before the learning criterion is reached, judgment error is comparable in the linear and the multiplicative judgment task in most conditions. Second, if participants remembered their judgments less well because they solved a fewer number of training blocks, participants who needed more training blocks to reach the learning criterion should show a lower rate of forgetting. Yet, participants who reached the learning criterion in the multiplicative task after 15 or more training blocks made on average more errors than participants who reached the criterion earlier in training. Furthermore, judgment error increased more strongly from training to test for participants who needed more blocks. In sum, those results make it unlikely that training stopped too early in the multiplicative task.

Implications for training

From a broader perspective, considering which knowledge people are more likely to forget may inform our understanding about how people can best acquire this knowledge and retain it for a long time. For instance, the knowledge about categories that people retain after a longer time interval depends on how they learned the task (Sakamoto & Love, 2010). Yet people do not always structure their learning in a way that facilitates later retrieval, neither in education (Bjork, Dunlosky, & Kornell, 2013) nor when learning abstract concepts (Tauber, Dunlosky, Rawson, Wahlheim, & Jacoby, 2013). Our study contributes to a new branch of research in function learning and categorization studying how to construct specific training procedures to improve categorization decisions over a

long time interval. This line of research has investigated how manipulations that improve long-term retention may help category or function learning and generalization, ranging from spaced training (Carvalho & Goldstone, 2014; McDaniel, Fadler, & Pashler, 2013; Zulkipli & Burt, 2013) to testing effects (Kang et al., 2011) to optimal training exemplars (Giguere & Love, 2013; Hornsby & Love, 2014). For instance, spacing exemplar presentations improves memory performance for trained items and simplifies generalization to new items (McDaniel et al., 2013). Our study emphasizes that identifying the underlying task structure and the strategies people use to approach the task can help to adapt those training procedures. Specifically, if rules can be abstracted as in linear judgment tasks, it may be sufficient to distribute training across time to achieve high judgment accuracy and adequate generalization. In contrast, multiplicative tasks require that participants identify and retrieve specific exemplars. Introducing new probes interferes with retrieval of those exemplars suggesting that successful training procedures potentially need to tackle this identification problem.

Conclusions

Since Ebbinghaus's (1885) seminal work much research has been devoted to the study of forgetting. Our study highlights that forgetting prior knowledge can similarly restrict how accurately people make judgments after some time has passed—not only if people need to retrieve past experiences, but also if they need to established a judgment policy based on abstracted knowledge. Identifying how abstracted knowledge and past experiences can best be retained may thus help improve human judgments in different domains from weather forecasts to business.

References

- Anderson, M. C., Bjork, R. A., & Bjork, E. L. (1994). Remembering can cause forgetting: Retrieval dynamics in long-term memory. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, *20*(5), 1063–1087.
doi:10.1037/0278-7393.20.5.1063
- Anderson, M. C. & Neely, J. H. (1996). Interference and inhibition in memory retrieval. In E. L. Bjork & R. A. Bjork (Eds.), *Memory* (2nd ed., pp. 237–313). San Diego, CA, US: Elsevier. doi:10.1016/B978-012102570-0/50010-0
- Ashton, R. H. (2000). A review and analysis of research on the test-retest reliability of professional judgment. *Journal of Behavioral Decision Making*, *13*(3), 277–294.
doi:10.1002/1099-0771(200007/09)13:3<277::AID-BDM350>3.0.CO;2-B
- Bahrick, H. P. (1984). Semantic memory content in permastore: Fifty years of memory for Spanish learned in school. *Journal of Experimental Psychology: General*, *113*(1), 1–29. doi:10.1037/0096-3445.113.1.30
- Bahrick, H. P., Bahrick, P. O., & Wittlinger, R. P. (1975). Fifty years of memory for names and faces: A cross-sectional approach. *Journal of Experimental Psychology: General*, *104*(1), 54–75. doi:10.1037/0096-3445.104.1.54
- Balzer, W. K., Rohrbaugh, J., & Murphy, K. R. (1983). Reliability of actual and predicted judgments across time. *Organizational Behavior & Human Performance*, *32*(1), 109–123. doi:10.1016/0030-5073(83)90142-3
- Bjork, R. A., Dunlosky, J., & Kornell, N. (2013). Self-regulated learning: Beliefs, techniques, and illusions. *Annual Review of Psychology*, *64*(1), 417–444.
doi:10.1146/annurev-psych-113011-143823
- Bourne, L. E., Healy, A. F., Kole, J. A., & Graham, S. M. (2006). Strategy shifts in classification skill acquisition: Does memory retrieval dominate rule use? *Memory & Cognition*, *34*(4), 903–913. doi:10.3758/BF03193436

- Brehmer, B. (1994). The psychology of linear judgement models. *Acta Psychologica*, *87*(2-3), 137–154. doi:10.1016/0001-6918(94)90048-5
- Brennan, R. L. & Prediger, D. J. (1981). Coefficient Kappa: Some uses, misuses, and alternatives. *Educational and Psychological Measurement*, *41*(3), 687–699. doi:10.1177/001316448104100307
- Capaldi, E. J. & Neath, I. (1995). Remembering and forgetting as context discrimination. *Learning & Memory*, *2*(3-4), 107–132. doi:10.1101/lm.2.3-4.107
- Carvalho, P. F. & Goldstone, R. L. (2014). Putting category learning in order: Category structure and temporal arrangement affect the benefit of interleaved over blocked study. *Memory & Cognition*, *42*(3), 481–495. doi:10.3758/s13421-013-0371-0
- Cooksey, R. W. (1996). *Judgment analysis: Theory, methods and applications*. San Diego, CA: Academic Press.
- Davis, T., Love, B. C., & Preston, A. R. (2012). Learning the exception to the rule: Model-based fMRI reveals specialized representations for surprising category members. *Cerebral Cortex*, *22*(2), 260–273. doi:10.1093/cercor/bhr036
- Ebbinghaus, H. (1885). *Über das Gedächtnis*. Leipzig: Duncker und Humblot.
- Einhorn, H. J., Kleinmuntz, D. N., & Kleinmuntz, B. (1979). Linear regression and process-tracing models of judgment. *Psychological Review*, *86*(5), 465–485. doi:10.1037//0033-295X.86.5.465
- Erickson, M. A. & Kruschke, J. K. (1998). Rules and exemplars in category learning. *Journal of Experimental Psychology: General*, *127*(2), 107–140. doi:10.1037/0096-3445.127.2.107
- Estes, W. K. (1986). Memory storage and retrieval processes in category learning. *Journal of Experimental Psychology: General*, *115*(2), 155–174. doi:10.1037/0096-3445.115.2.155
- French, R. M. (1999). Catastrophic forgetting in connectionist networks. *Trends in Cognitive Sciences*, *3*(4), 128–135. doi:10.1016/S1364-6613(99)01294-2

- Giguere, G. & Love, B. C. (2013). Limits in decision making arise from limits in memory retrieval. *Proceedings of the National Academy of Sciences*, *110*(19), 7613–7618.
doi:10.1073/pnas.1219674110
- Hahn, U. & Chater, N. (1998). Similarity and rules: Distinct? Exhaustive? Empirically distinguishable? *Cognition*, *65*(2-3), 197–230. doi:10.1016/S0010-0277(97)00044-9
- Hoffmann, J. A., von Helversen, B., & Rieskamp, J. (2014). Pillars of judgment: How memory abilities affect performance in rule-based and exemplar-based judgments. *Journal of Experimental Psychology: General*, *143*(6), 2242–2261.
doi:10.1037/a0037989
- Hoffmann, J. A., von Helversen, B., & Rieskamp, J. (2016). Similar task features shape judgment and categorization processes. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, *42*(8), 1193–1217. doi:10.1037/xlm0000241
- Homa, D., Cross, J., Cornell, D., Goldman, D., & Shwartz, S. (1973). Prototype abstraction and classification of new instances as a function of number of instances defining the prototype. *Journal of Experimental Psychology*, *101*(1), 116–122.
doi:10.1037/h0035772
- Hornsby, A. N. & Love, B. C. (2014). Improved classification of mammograms following idealized training. *Journal of Applied Research in Memory and Cognition*, *3*(2), 72–76. doi:10.1016/j.jarmac.2014.04.009
- Janssen, D. P. (2012). Twice random, once mixed: Applying mixed models to simultaneously analyze random effects of language and participants. *Behavior Research Methods*, *44*(February), 232–247. doi:10.3758/s13428-011-0145-1
- Juslin, P., Karlsson, L., & Olsson, H. (2008). Information integration in multiple cue judgment: A division of labor hypothesis. *Cognition*, *106*(1), 259–298.
doi:10.1016/j.cognition.2007.02.003

- Juslin, P., Olsson, H., & Olsson, A.-C. (2003). Exemplar effects in categorization and multiple-cue judgment. *Journal of Experimental Psychology: General*, *132*(1), 133–156. doi:10.1037/0096-3445.132.1.133
- Kang, S. H. K., McDaniel, M. A. M., & Pashler, H. (2011). Effects of testing on learning of functions. *Psychonomic Bulletin and Review*, *18*(5), 998–1005. doi:10.3758/s13423-011-0113-x
- Karlsson, L., Juslin, P., & Olsson, H. (2007). Adaptive changes between cue abstraction and exemplar memory in a multiple-cue judgment task with continuous cues. *Psychonomic Bulletin & Review*, *14*(6), 1140–1146. doi:10.3758/BF03193103
- Karpicke, J. D. & Roediger, H. L. (2008). The critical importance of retrieval for learning. *Science*, *319*(5865), 966–968. doi:10.1126/science.1152408
- Kruschke, J. K. (1992). ALCOVE: An exemplar-based connectionist model of category learning. *Psychological Review*, *99*(1), 22–44. doi:10.1037/0033-295X.99.1.22
- Little, J. L. & McDaniel, M. A. M. (2015). Individual differences in category learning: Memorization versus rule abstraction. *Memory & Cognition*, *43*(2), 283–297. doi:10.3758/s13421-014-0475-1
- Mata, R., von Helversen, B., Karlsson, L., & Cüpper, L. (2012). Adult age differences in categorization and multiple-cue judgment. *Developmental Psychology*, *48*(4), 1188–1201. doi:10.1037/a0026084
- McDaniel, M. A. M., Fadler, C. L. C., & Pashler, H. (2013). Effects of spaced versus massed training in function learning. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, *39*(5), 1417–1432. doi:10.1037/a0032184
- McGeoch, J. A. (1932). Forgetting and the law of disuse. *Psychological Review*, *39*(4), 352–370. doi:10.1037/h0069819
- Meeter, M., Murre, J. M. J., & Janssen, S. M. J. (2005). Remembering the news: Modeling retention data from a study with 14,000 participants. *Memory & Cognition*, *33*(5), 793–810. doi:10.3758/BF03193075

- Morris, S. B. & DeShon, R. P. (2002). Combining effect size estimates in meta-analysis with repeated measures and independent-groups designs. *Psychological Methods*, 7(1), 105–125. doi:10.1037/1082-989X.7.1.105
- Nosofsky, R. M. (1988). Exemplar-based accounts of relations between classification, recognition, and typicality. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 14(4), 700–708. doi:10.1037//0278-7393.14.4.700
- Nosofsky, R. M. & Johansen, M. K. (2000). Exemplar-based accounts of "multiple-system" phenomena in perceptual categorization. *Psychonomic Bulletin & Review*, 7(3), 375–402.
- Nosofsky, R. M. & Zaki, S. R. (1998). Dissociations between categorization and recognition in amnesic and normal individuals: An exemplar-based interpretation. *Psychological Science*, 9(4), 247–255. doi:10.1111/1467-9280.00051
- Olsson, A.-C., Enkvist, T., & Juslin, P. (2006). Go with the flow: How to master a nonlinear multiple-cue judgment task. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 32(6), 1371–1384. doi:10.1037/0278-7393.32.6.1371
- Pachur, T. & Olsson, H. (2012). Type of learning task impacts performance and strategy selection in decision making. *Cognitive Psychology*, 65(2), 207–240. doi:10.1016/j.cogpsych.2012.03.003
- Palmeri, T. J. & Nosofsky, R. M. (1995). Recognition memory for exceptions to the category rule. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 21(3), 548–568. doi:10.1037/0278-7393.21.3.548
- Platzer, C. & Bröder, A. (2013). When the rule is ruled out: Exemplars and rules in decisions from memory. *Journal of Behavioral Decision Making*, 26(5), 429–441. doi:10.1002/bdm.1776
- Posner, M. I. & Keele, S. W. (1970). Retention of abstract ideas. *Journal of Experimental Psychology*, 83(2), 304–308.

- Postman, L. (1971). Organization and interference. *Psychological Review*, *78*(4), 290–302.
doi:10.1037/h0031031
- Richardson, D. C. & Spivey, M. J. (2000). Representation, space and Hollywood Squares: Looking at things that aren't there anymore. *Cognition*, *76*(3), 269–295.
doi:10.1016/S0010-0277(00)00084-6
- Robbins, D., Barresi, J., Compton, P., Furst, A., Russo, M., & Smith, M. A. (1978). The genesis and use of exemplar vs. prototype knowledge in abstract category learning. *Memory & Cognition*, *6*(4), 473–480. doi:10.3758/BF03197481
- Roediger, H. L. & Karpicke, J. D. (2006). The power of testing memory. *Perspectives on Psychological Science*, *1*(3), 181–210. doi:10.1111/j.1745-6916.2006.00012.x
- Roediger, H. L., Weinstein, Y., & Agarwal, P. K. (2010). Forgetting. In S. Della Sala (Ed.), *Forgetting: Preliminary considerations* (pp. 1–22). Hove, UK: Psychology Press.
doi:10.4324/9780203851647
- Rouder, J. N. & Ratcliff, R. (2004). Comparing categorization models. *Journal of Experimental Psychology: General*, *133*(1), 63–82. doi:10.1037/0096-3445.133.1.63
- Rubin, D. C. & Wenzel, A. E. (1996). One hundred years of forgetting: A quantitative description of retention. *Psychological Review*, *103*(4), 734–760.
doi:10.1037/0033-295X.103.4.734
- Sakamoto, Y. & Love, B. C. (2004). Schematic influences on category learning and recognition memory. *Journal of Experimental Psychology: General*, *133*(4), 534–553.
doi:10.1037/0096-3445.133.4.534
- Sakamoto, Y. & Love, B. C. (2010). Learning and retention through predictive inference and classification. *Journal of Experimental Psychology: Applied*, *16*(4), 361–377.
doi:10.1037/a0021610
- Scholz, A., von Helversen, B., & Rieskamp, J. (2015). Eye movements reveal memory processes during similarity- and rule-based decision making. *Cognition*, *136*, 228–246.
doi:10.1016/j.cognition.2014.11.019

- Schwarz, G. (1978). Estimating the dimension of a model. *The Annals of Statistics*, 6(2), 461–464.
- Siegel, J. A. & Siegel, W. (1972). Absolute judgment and paired-associate learning: Kissing cousins or identical twins? *Psychological Review*, 79(4), 300–316.
doi:10.1037/h0032945
- Tauber, S. K., Dunlosky, J., Rawson, K. a., Wahlheim, C. N., & Jacoby, L. L. (2013). Self-regulated learning of a natural category: Do people interleave or block exemplars during study? *Psychonomic Bulletin & Review*, 20(2), 356–63.
doi:10.3758/s13423-012-0319-6
- von Helversen, B., Mata, R., & Olsson, H. (2010). Do children profit from looking beyond looks? From similarity-based to cue abstraction processes in multiple-cue judgment. *Developmental Psychology*, 46(1), 867–889. doi:10.1037/a0016690
- von Helversen, B. & Rieskamp, J. (2008). The mapping model: A cognitive theory of quantitative estimation. *Journal of Experimental Psychology: General*, 137(1), 73–96.
doi:10.1037/0096-3445.137.1.73
- von Helversen, B. & Rieskamp, J. (2009). Models of quantitative estimations: Rule-based and exemplar-based processes compared. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 35(4), 867–889. doi:10.1037/a0015501
- Wagenmakers, E.-J. & Farrell, S. (2004). AIC model selection using Akaike weights. *Psychonomic Bulletin & Review*, 11(1), 192–196. doi:10.3758/BF03206482
- Zulkipli, N. & Burt, J. S. (2013). The exemplar interleaving effect in inductive learning: Moderation by the difficulty of category discriminations. *Memory & Cognition*, 41(1), 16–27. doi:10.3758/s13421-012-0238-9

Table 1

Task Structure in Experiment 1

Cue values				Criterion		Item Type
C_1	C_2	C_3	C_4	Linear	Multiplicative	
0	1	1	0	15	12	Old
0	0	0	1	11	10	Old
1	1	1	0	19	18	Old
0	0	1	1	13	11	Old
0	1	0	1	14	12	Old
0	0	0	0	10	10	Old
1	1	1	1	20	20	Old
1	0	1	0	16	13	Old
0	1	1	1	16	13	Old
0	0	1	0	12	11	Old
1	0	1	1	17	14	New
1	1	0	0	17	14	New
0	1	0	0	13	11	New
1	1	0	1	18	16	New
1	0	0	1	15	12	New
1	0	0	0	14	12	New

Note. The judgment criterion was derived from Equation 1 (linear) and Equation 2 (multiplicative).

Table 2

Performance in Experiment 2. Standard Error in Parentheses.

	Judgment Task			
	Linear		Multiplicative	
	Retention interval		Retention interval	
	Immediate	1 week	Immediate	1 week
Training session				
Number of blocks	15.9 (0.7)	16.9 (0.6)	14.6 (0.7)	15.7 (0.7)
Error last block	1.33 (0.16)	1.52 (0.21)	1.34 (0.21)	1.63 (0.23)
Test session				
Error training items	1.24 (0.16)	1.71 (0.18)	1.49 (0.16)	2.37 (0.18)
Error validation items	1.86 (0.19)	2.01 (0.17)	2.42 (0.14)	2.79 (0.16)
Recognition (% correct)	61.6 (3.4)	69.2 (2.9)	62.6 (3.1)	58.4 (2.9)

Note. Error was measured as the RMSD (Root Mean Squared Deviation) between participant's judgment and the criterion.

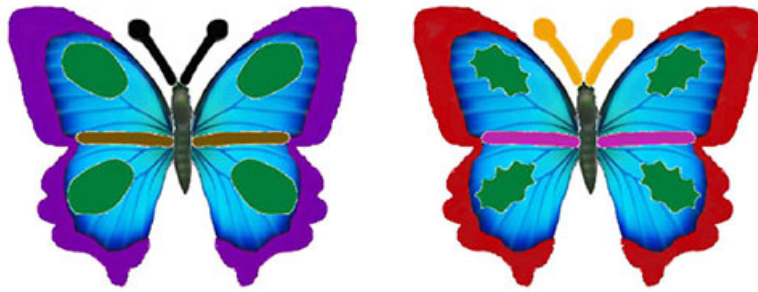


Figure 1. Sample species of butterflies with distinct cue values on all cues.

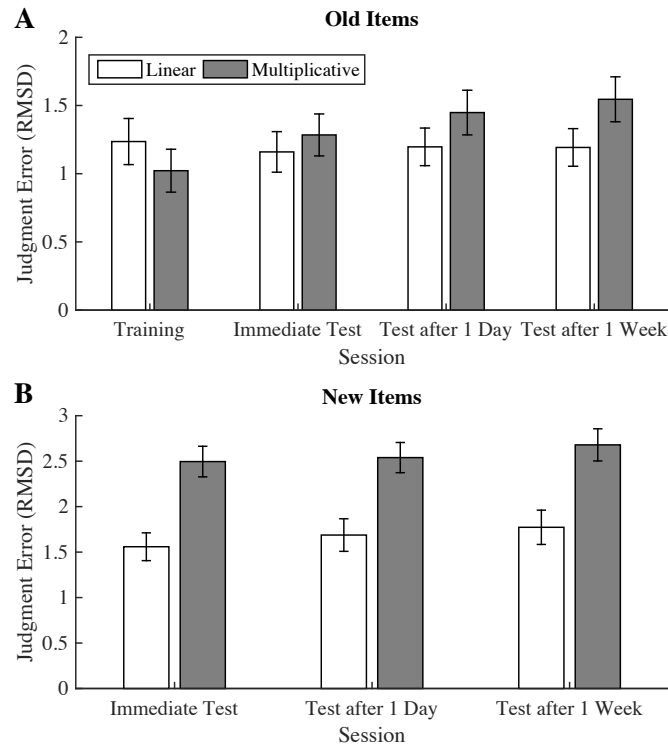


Figure 2. Judgment error measured in root-mean squared deviation (RMSD) on old items (Panel A) and new items (Panel B) in Experiment 1. White bars depict judgment error in the linear judgment task, gray bars depict judgment error in the multiplicative judgment task. **A.** Judgment error on old items was assessed for each participant in the last block of training as well as in all three test sessions (immediate test, test after 1 day, test after 1 week). **B.** Judgment error on new items was assessed for each participant in all three test sessions. Error bars indicate $\pm 2 SE$.

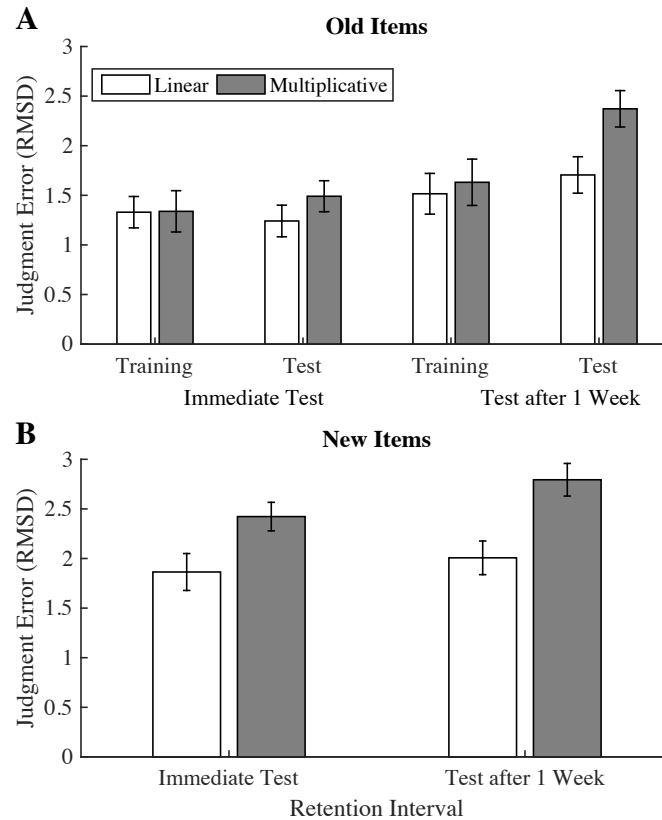


Figure 3. Judgment error measured in root-mean squared deviation (RMSD) on old items (Panel A) and new items (Panel B) in Experiment 2. White bars depict judgment error in the linear judgment task, gray bars depict judgment error in the multiplicative judgment task. **A.** Judgment error on old items was assessed for each participant in the last block of training as well as after a short retention interval (immediate test) or a long retention interval (test after one week). **B.** Judgment error on new items was assessed for each participant after either a short retention interval (immediate test) or a long retention interval (test after one week). Error bars indicate $\pm 2 SE$

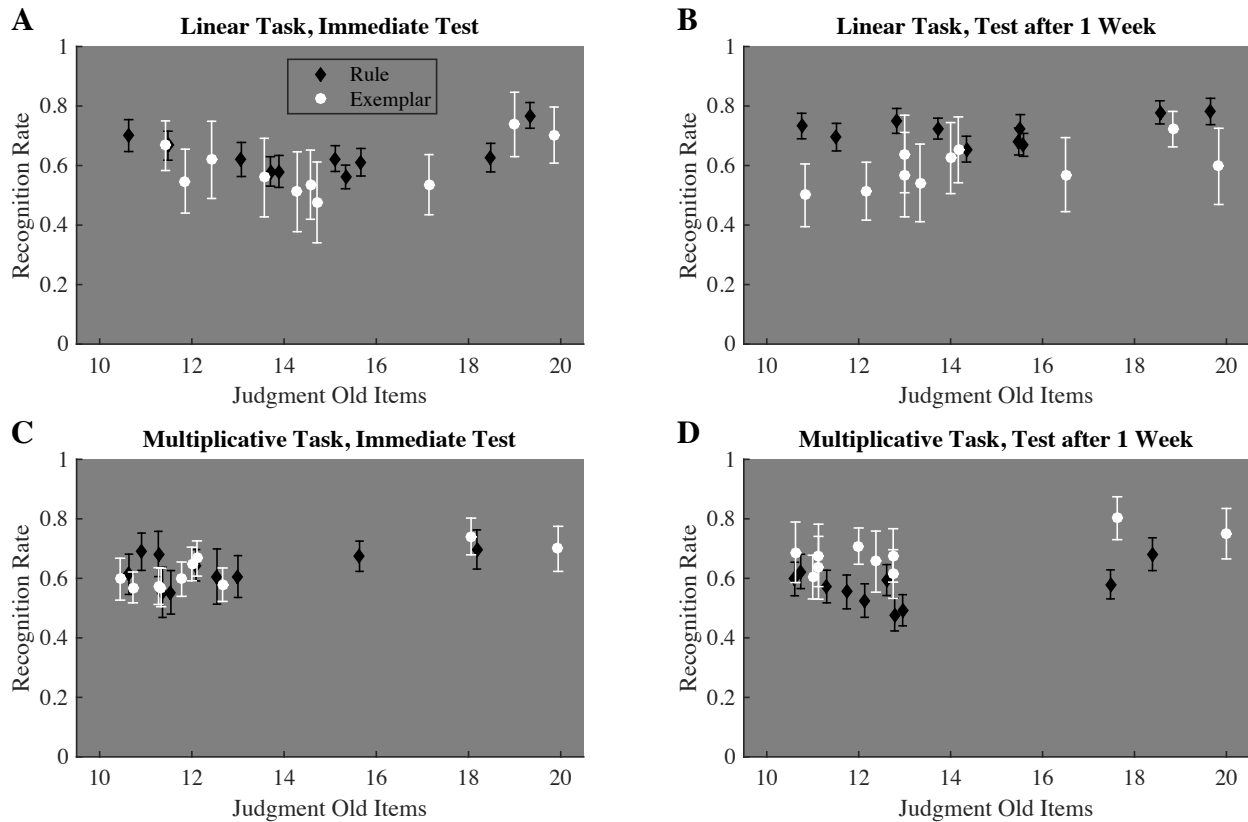


Figure 4. Proportion of correctly recognized old items (recognition rate) plotted against the average judgment for this old item in the last block of training, separately for participants classified to the rule-based (black diamonds) and the exemplar-based judgment strategy (white circles). Panel **A** depicts the recognition rate for participants who solved the linear task and took an immediate test. Panel **B** depicts the recognition rate for participants who solved the linear task and took the test after a week. Panel **C** and **D** show the recognition rate in the multiplicative task for immediate test and test after a week, respectively. Error bars indicate $\pm 2 SE$.

Appendix A

Cognitive modeling of judgment strategies

To identify the cognitive strategies that people rely on in the three test sessions, we used a computational modeling approach. We compared how well a prominent rule-based model, the linear model, described participants' judgments in comparison to one often-used exemplar model using four attention weights. We compared all models to a guessing model that assumed that participants' judgments vary around participants' mean judgment on each trial. The guessing model estimated two free parameters: participants' mean judgment and the fitted standard deviation.

Model description

Rule-based model. To model rule-based strategies, we fitted a linear, additive model that has often served as the prototypical rule-based strategy in judgment tasks (Cooksey, 1996; Juslin et al., 2003). Corresponding mathematically to a linear regression, the linear model allows combining several cues in a linear additive fashion. Accordingly, the estimated criterion value \hat{c}_p of an object p is the weighted sum of the cue values x_{pi} :

$$\hat{c}_p = k + \sum_{i=1}^4 w_i \cdot x_{pi} \quad (3)$$

where w_i are the cue weights for each cue i and k is a constant intercept. In sum, the linear model estimates six parameters: four cue weights w_i , one intercept k , and the standard deviation.

Exemplar model. Exemplar models have been widely used in judgment and categorization research to model retrieval of single instances from long-term memory (Hoffmann et al., 2014; Juslin et al., 2003). In exemplar models, the similarity $S(p, j)$ between the probe p and exemplar j is an exponential decay function of the distances d_{pj} between the objects (Nosofsky & Zaki, 1998).

$$S(p, j) = e^{-d_{pj}} \quad (4)$$

Thus, smaller distances between the probe p and exemplar j indicate a higher similarity between these objects. To determine this distance, the cue values x_{pi} of probe p are compared to the cue values x_{ji} of exemplar j on all cues i . The more the cue values match each other, the smaller is the distance between the objects (Nosofsky & Johansen, 2000).

$$d_{pj} = h \left(\sum_{i=1}^4 w_i |x_{pi} - x_{pj}| \right) \quad (5)$$

The sensitivity parameter h determines how strongly similarity decays with distance. Smaller sensitivity parameters indicate that similarity declines less with distance. The attention weights w_i , summing to one, weigh how much attention each cue or dimension receives. To account for judgments, Juslin et al. (2003) assumed that the criterion value c_j of an exemplar is stored together with its cue values in memory. To estimate the criterion value of a new probe \hat{c}_p , the criterion values c_j for each exemplar are weighted by the similarities.

$$\hat{c}_p = \frac{\sum_{j=1}^J S(p, j) \cdot c_j}{\sum_{j=1}^J S(p, j)} \quad (6)$$

In sum, the exemplar model estimates five parameters: three attention weights w_i , the sensitivity parameter h , and the standard deviation.

Model estimation and comparison. To evaluate the models' relative performance we fitted all models to participants' judgment on all six presentations of old and new items, separately for each of the three test sessions (immediate test, test after a day, and test after a week). The models were evaluated based upon the Bayesian Information Criterion (BIC; Schwarz (1978)). All models were fitted to participants' responses by minimizing the deviance $-2LL$, the negative summed log-likelihood L of the model given the data.

$$-2LL = -2 \cdot \sum \ln(L) \quad (7)$$

We calculated the likelihood as the probability density of participants' judgments j assuming a truncated normal distribution, with the models' predicted responses \hat{c}_p as the

mean of the normal distribution and a fitted standard deviation σ . This truncated normal distribution was chosen because it matched the response scale from 10 to 20.

$$L = \frac{1}{\sigma} \frac{\phi(j|\hat{c}_p, \sigma)}{\Phi(20|\hat{c}_p, \sigma) - \Phi(10|\hat{c}_p, \sigma)} \quad (8)$$

To compare which model described participants' responses better, we calculated the BIC for each model. The BIC can be used to compare non-nested models and penalizes more complex models by accounting for the number of free model parameters k :

$$\text{BIC} = -2LL + k \ln n, \quad (9)$$

where n denotes the number of observations. Smaller BIC values indicate a better model fit. BICs were converted into BIC weights $\text{BIC}_{w,M}$ that give the posterior probability of each model given the data (Wagenmakers & Farrell, 2004).

$$\text{BIC}_{w,M} = \frac{e^{-.5\Delta\text{BIC}_M}}{\sum_i e^{-.5\Delta\text{BIC}_i}} \quad (10)$$

with ΔBIC_M as the difference between model M and the best model in the set and ΔBIC_i as the difference between a specific model i the best model.

Detailed results for model comparisons in Experiment 1

Table A1 shows BIC weights, strategy classifications based on the BIC weights as well as the RMSD between model responses and participants' judgments. Across all sessions and judgment tasks, only a few participants were best described by a guessing model, as low BIC weights and a low number of participants classified to the guessing model indicate. Furthermore, BIC weights were higher for the linear model than for the exemplar model in all sessions and in both the linear and the multiplicative judgment task and more participants were classified to the linear model than to the exemplar model. Compared to the linear task, however, BIC weights were higher for the exemplar model in the multiplicative judgment task and more participants were best described by an exemplar model. Furthermore, the exemplar model provided a better fit to participants'

judgments in the multiplicative judgment task, whereas the linear model fitted participants' judgments better in the linear task.

Table A1

Strategy classification in Experiment 1. SE for BIC weights and RMSD in Parantheses

Task	Session	BIC _w			Classification			RMSD Model response			
		Guess	Rule	Ex	Guess	Rule	Ex	Guess	Rule	Ex	Weighted
Lin	Immediate	.05 (.03)	.79 (.06)	.16 (.06)	2	32	6	2.48 (0.08)	0.65 (0.06)	0.80 (0.04)	0.63 (0.06)
	One day	.05 (.03)	.74 (.07)	.21 (.06)	2	29	9	2.53 (0.09)	0.63 (0.07)	0.82 (0.05)	0.59 (0.07)
	One week	.02 (.02)	.85 (.06)	.12 (.05)	1	34	5	2.62 (.08)	0.62 (0.07)	0.86 (0.05)	0.60 (0.07)
Mult	Immediate	.03 (.02)	.59 (.08)	.39 (.08)	1	24	15	2.75 (0.09)	1.48 (0.08)	1.07 (0.09)	1.12 (0.09)
	One day	.02 (.02)	.58 (.08)	.40 (.08)	1	23	16	2.85 (0.08)	1.56 (0.17)	1.07 (0.10)	1.08 (0.16)
	One week	.00 (.00)	.65 (.07)	.35 (.08)	0	26	14	2.88 (0.08)	1.49 (0.12)	1.12 (0.09)	1.18 (0.13)

Note. BIC_w = Bayesian Information Criterion weights, RMSD = Root Mean Squared Deviation, Lin = Linear task, Mult = Multiplicative task, Guess = Guessing model, Rule = Linear model, Ex = Exemplar model. Weighted RMSD indicates the RMSD between participants' judgments and a model that weighs model responses of the guessing, linear, and exemplar model with their corresponding BIC_w.

Detailed results for model comparisons in Experiment 2

Table A2 displays BIC weights, strategy classifications based on BIC weights, and the RMSD between model responses and participants' judgments. Similar to Experiment 1, BIC weights for the guessing model as well as the number of participants classified to the guessing model were low with a slightly higher number of participants classified to the guessing model in the multiplicative task. In the linear task, the majority of participants were best described by the linear model both in the immediate test session as well as in test after one week, as shown by high average BIC weights for the linear model and a high number of participants classified to that model. The exemplar model better described participants' judgment in the immediate test in the multiplicative task, as suggested by higher BIC weights and more participants classified to the model. In test after one week, however, the linear model provided a higher BIC weight and more participants were classified to the linear model.

Table A2

Strategy classification in Experiment 2. Standard Error for BIC_w and RMSD in Parantheses

Task	Session	BIC_w			Classification			RMSD Model response			
		Guess	Rule	Ex	Guess	Rule	Ex	Guess	Rule	Ex	Weighted
Lin	Immediate	.03 (.03)	.79 (.07)	.19 (.06)	1	27	7	2.58 (0.08)	0.69 (0.07)	0.88 (0.07)	0.67 (0.07)
	One week	.03 (.03)	.82 (.06)	.16 (.06)	1	29	6	2.51 (0.10)	0.74 (0.07)	1.08 (0.09)	0.74 (0.07)
Mult	Immediate	.07 (.04)	.37 (.08)	.56 (.08)	2	11	18	2.56 (0.14)	1.83 (0.16)	0.97 (0.07)	1.06 (0.11)
	One week	.10 (.05)	.65 (.08)	.25 (.07)	4	23	8	2.35 (0.15)	1.42 (0.13)	1.45 (0.11)	1.24 (0.13)

Note. BIC_w = Bayesian Information Criterion weights, RMSD = Root Mean Squared Deviation, Lin = Linear task, Mult = Multiplicative task, Guess = Guessing model, Rule = Linear model, Ex = Exemplar model. Weighted RMSD indicates the RMSD between participants' judgments and a model that weighs model responses of the guessing, linear, and exemplar model with their corresponding BIC_w .

Appendix B

Modeling forgetting of exemplars

We modeled forgetting of exemplars by using a successful model that learns to make exemplar-based decisions (Kruschke, 1992). The model assumes that learning in judgment tasks can be understood as gradually forming associative links between the exemplars that are encountered and the possible criterion values. Judgments are a function of the similarity of the probe to the previously encountered exemplars and of the association strengths between the exemplars and the criterion values. That is, the probe activates similar exemplars, which in turn activate criterion values they are associated with.

Association strengths are then translated into output probabilities for each criterion value and the final judgment is the mean of the criterion values weighted by their probabilities.

The model contains three free parameters, two learning parameters and a sensitivity parameter: The first learning parameter determines the speed with which the associations between criterion values and exemplars are formed. The second learning parameter determines how fast people learn to differentially distribute attention to the features of the objects and changes how similarity between the probe and learnt exemplars is computed. The sensitivity parameter regulates how similarity is translated into the activation of an exemplar.

We introduced forgetting in this exemplar-based learning model by assuming that over time further exemplars would be encountered that interfere with previously learnt information. New information updates both the association between exemplars and stored criterion values as well as learned attention towards specific cues. For each simulation, we randomly drew 1000 times from an exponential distribution for the two learning parameters (with $M = 1$) and the sensitivity parameter (with $M = 3$). Training followed the same schedule as in the experiment, but we introduced four additional random cues in the simulation to limit catastrophic forgetting (French, 1999), the effect that learning new information completely erases previously learned information. Between training and test,

the exemplar model continued to learn N random item profiles (N varied from 0 to 100 in steps of 10). Finally, the model made the same judgments for old exemplars in test as participants did.