# Group beats Trend!?
# Testing feature hierarchy in statistical graphics

Susan VanderPlas, Heike Hofmann[*]

March 5, 2015

**Abstract**

abstract goes here

## 1   Introduction and background

Discussion of pre-attentive visual features (Healey and Enns, 2012) - with a focus on hierarchy of pre-attentive features: color trumps shape - do we also see this in our results, and if so, by how much?

Numerical information can be difficult to communicate effectively in raw form, due to limits on attention span, short term memory, and information storage mechanisms within the human brain. Graphics are much more effective for communicating numerical information, as (well-designed) graphics order the numerical information spatially and utilize the higher-bandwidth visual system. Visual data displays serve as a form of external cognition **??**, ordering and visually summarizing data which would be hopelessly confusing in tabular format. One fantastic example of this phenomenon is the Hertzsprung-Russell (HR) diagram, which was described as "one of the greatest observational syntheses in astronomy and astrophysics" because it allowed astronomers to clearly relate the absolute magnitude of a star to its' spectral classification; facilitating greater understanding of stellar evolution (Spence and Garrison, 1993). The data it displayed was previously available in several different tables; when plotted on the same chart, information that was invisible in a tabular representation became immediately clear (Lewandowsky and Spence, 1989b). Graphical displays more efficiently utilize cognitive resources by reducing the burden of storing, ordering, and summarizing raw data; this frees bandwidth for higher levels of information synthesis, allowing observers to note outliers, understand relationships between variables, and form new hypotheses.

Graphical displays are powerful because they efficiently and effectively convey numerical information, but we have relatively sparse empirical information about how the human perceptual system processes these displays. Our understanding of the perception of statistical graphics is informed by general psychological and psychophysics research as well as more specific research into the perception of data displays(Cleveland and McGill, 1984).

---

[*]Department of Statistics and Statistical Laboratory, Iowa State University

One relevant focus of psychological research is pre-attentive perception, that is, perception which occurs automatically in the first 200 ms of exposure to a visual stimulus (Treisman, 1985).

Research into preattentive perception provides us with some information about the temporal hierarchy of graphical feature processing. Color, line orientation, and shape are processed preattentively; that is, within 200 ms, it is possible to identify a single target in a field of distractors, if the target differs with respect to color or shape (Goldstein, 2009). Research by Healey and Enns (1999) extends this work, demonstrating that certain features of three-dimensional data displays are processed preattentively, but that neither target identification nor three-dimensional data processing always translate into faster or more accurate inference about the data displayed, particularly when participants must integrate several preattentive features to understand the data.

Feature detection at the attentive stage of perception has also been examined in the context of statistical graphics; researchers have evaluated the perceptual implications of utilizing color, fill, shapes, and letters to denote categorical or stratified data in scatterplots. Cleveland and McGill (1984) ranked the optimality of these plot aesthetics based on response accuracy, preferring colors, amount of fill, shapes, and finally letters to indicate category membership. Lewandowsky and Spence (1989a) examined both accuracy and response time, finding that color is faster and more accurately perceived (except for those with color deficiency). Shape, fill, and discriminable letters (letters which do not share visual features, such as HQX) were identified as less accurate than color, while confusable letters (such as HEF) result in significantly decreased accuracy.

Another area of psychological research, Gestalt psychology, examine perception as a holistic experience, establishing and evaluating mental heuristics used to transform visual stimuli into useful, coherent information. Gestalt rules of perception can be easily applied to statistical graphics, as they describe the way we organize visual input, focusing on the holistic experience rather than the individual perceptual features.

For example, rather than perceiving four legs, a tail, two eyes, two ears, and a nose, we perceive a dog. The rules of perceptual grouping or organization, as stated in Goldstein (2009) are:

- Proximity: two elements which are close together are more likely to belong to a single unit.

- Similarity: the more similar two elements are, the more likely they belong to a single unit.

- Common fate: two elements moving together likely belong to a single unit.

- Good continuation: two elements which blend together smoothly likely belong to one unit.

- Closure: elements which can be assembled into closed or convex objects likely belong together.

- Common region: elements contained within a common region likely belong together.

- Connectedness: elements physically connected to each other are more likely to belong together.

The plots in figure 1 demonstrate several of the gestalt principles which combine to order our perceptual experience from the top down. These laws help to order our perception of charts as well: points which are colored or shaped the same are perceived as belonging to a group (similarity), points within a bounding interval or ellipse are perceived as belonging to the same group (common
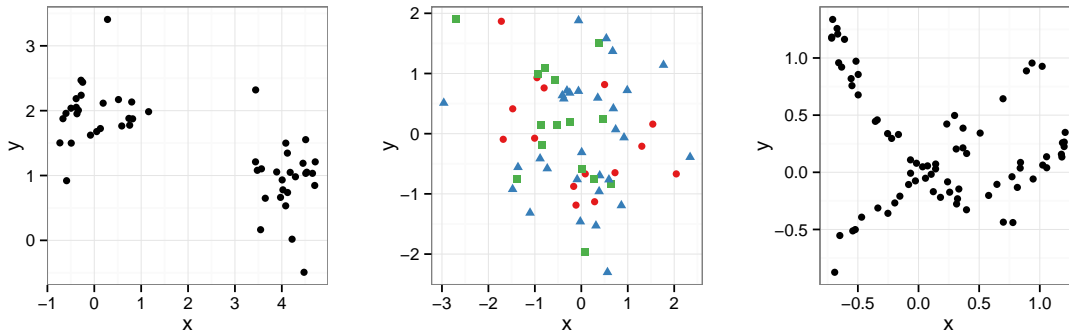
Figure 1: *Proximity* renders the fifty points of the first scatterplot as two distinct (and equal-sized) groups. Shapes and colors create different groups of points in the middle scatterplot, invoking the Gestalt principle of *Similarity*. *Good Continuation* renders the points in the scatterplot on the right hand side into two groups of points on curves: one a straight line with an upward slope, the other a curve that initially decreases and at the end of the range shows an uptick.

region), and regression lines with confidence intervals are perceived as single units (connectedness, closure, and/or common region).

clarify next sentence

The use of physical location, color, and shape to organize graphical units mentally utilizes both preattentive processing and higher-order gestalt schemas, identifying and grouping similar graphical features and simultaneously directing attention to graphical features which stand alone.

Research on preattentive perception is important because features that are perceived preattentively do not require as much mental effort to process from raw visual stimuli; theoretically, subsequent top-down gestalt heuristics can be applied to such stimuli more quickly.

We should also look at the time to response – it would be interesting, to see if the conflicting stimuli need more time to come to a decision. It's obviously not milliseconds that we measure, but it might still be informative. (We would need to exclude everybody's first attempt).

This study is designed to understand the hierarchy of gestalt principles in perception of statistical graphics. We utilize information from previous studies (Çağatay Demiralp et al., 2014; Robinson, 2003) concerning the hierarchy of preattentive feature perception in order to maximize the effect of preattentive feature differences.

might be useful to have a small diagram describing the perceptual process (with preattentive processing way at the top and gestalt heuristic processing in the middle, with "cognitive effort" at the bottom). Not sure if it's necessary, though. HH: good idea, let's see how much space we'll have.

Statistical graphics can be difficult to examine experimentally; qualitative studies rely on descriptions of the plot by participants who may not be able to articulate their inferences precisely, while quantitative studies may only be able to examine whether the viewer can accurately read numerical information from the chart, instead of exploring the overall utility of the data display wholistically. Statistical lineups, described in the next section, are an important experimental tool for evaluating the perceptual utility of graphical displays. Lineups fuse commonly used psychological tests (target identification, visual search) with statistical hypothesis tests to facilitate formal experimental evaluation of statistical graphics.

## 1.1 Statistical Lineups

Intro to lineups (**?**Majumder et al., 2013; **?**; **?**).

Describe the lineup protocol, including basic statistics. Link to the psychological "target and distractors" approach, which can be used to justify the addition of a second target, even with the PITA of the statistical complications.

In this study, we modify the lineup protocol by introducing a second target to each lineup. The two targets represent two different, competing signals; the participant's choice then demonstrates empirically which signal is more salient.

If both targets exhibit similar signal, participants may identify both targets, removing any forced-choice scenario which might skew results (few participants exercise this option).

By tracking the proportion of observers choosing either target plot (a measure of overall lineup difficulty) as well as which proportion of observers choose one target over the other target, we can determine the relative strength of the two competing signals amid a field of distractors. At this level, signal strength is determined by the experimental data and the generating model; we are measuring the "power" (in a statistical sense) of the human perceptual system, rather than raw numerical signal.

Using this testing framework, we can apply different aesthetics, such as color and shape, as well as plot objects which display statistical calculations, such as trend lines and bounding ellipses. These additional plot layers, discussed in more detail in the next section, are designed to emphasize one of the two competing targets and affect the overall visual signal of the target plot relative to the null plots. We expect that in a situation similar to the third plot of figure 1, the addition of two trend lines would emphasize the "good continuation" of points in the plot, producing a stronger visual signal, even though the underlying data has not changed. Similarly, the grouping effect in the first plot in the figure would be enhanced if the points in each group were colored differently, as the proximity heuristic would be supplemented by similarity. In plots that are ambiguous, containing some clustering of points as well as a linear relationship between $x$ and $y$, additional aesthetic cues may "tip the balance" in favor of recognizing one type of signal.

beautiful!!

This study is designed to inform our understanding of the perceptual implications of these additional aesthetics, in order to provide guidelines for the creation of data displays which provide visual cues consistent with gestalt heuristics and preattentive perceptual preferences. The next section discusses the particulars of the experimental design, including the data generation model, plot aesthetics, selection of color and shape palettes, and other important considerations.

# 2  Experimental Design

In this section, we discuss the generating data models for the two types of signal plots and the null plots, the selection of plot aesthetic combinations and aesthetic values, and the design and execution of the experiment.

> I know this will have to be rearranged, expanded, and transitions between sections will need to be added, but I want to get the paragraphs out.

## 2.1  Data Generation

Lineups require a single "target" data set (which we are expanding to two competing "target" data sets), and a method for generating null plots. When utilizing real data for target plots, null plots are often generated through bootstrap sampling, but this introduces some dependencies between target and null plots which complicate the statistical analysis of the results.

> add citations

When possible, it is desireable to generate true null plots, which are generated from the null model and do not depend on the data used in the target plot. This experiment will measure two competing gestalt heuristics, proximity and good continuation, using two data-generating models: $M_C$, which generates data with $K$ clusters, and $M_T$, which generates data with a positive correlation between $x$ and $y$. True null datasets are created using a mixture model $M_0$ which combines $M_C$ and $M_T$. Both $M_C$ and $M_T$ generate data in the same range of values. Additionally, $M_C$ generates clustered data with linear correlations that are within $\rho = (0.25, 0.75)$, similar to the linear relationship between datasets generated by $M_0$, and $M_T$ generates data with clustering similar to $M_0$. These constraints provide some assurance that participants who select a plot with data generated from $M_T$ are doing so because of visual cues indicating a linear trend (rather than a lack of clustering compared to plots with data generated from $M_0$), and participants who select a plot with data generated from $M_C$ are doing so because of visual cues indicating clustering, rather than a lack of a linear relationship relative to plots with data generated from $M_0$.

### 2.1.1  Regression Model $M_T$

This model has the parameter $\sigma_T$ to reflect the amount of scatter around the trend line. It generates $N$ points $(x_i, y_i)$, $i = 1, ..., N$ where $x$ and $y$ have a positive linear relationship. The data generation mechanism is as follows:

**Algorithm 2.1**
*Input Parameters: sample size $N$, $\sigma_T$ standard deviation around the line*
*Output: $N$ points, in form of vectors $x$ and $y$.*

1. *Generate $\tilde{x}_i$, $i = 1, ..., N$, as a sequence of evenly spaced points from $[-1, 1]$*

2. *Jitter $\tilde{x}_i$ by adding small uniformly distributed perturbations to each of the values: $x_i = \tilde{x}_i + \eta_i$, $\eta_i \sim Unif(-z, z)$, $z = 1/5 \cdot 2/(N-1)$*

3. *Generate $y_i$: $y_i = x_i + e_i$, $e_i \sim N(0, \sigma_T^2)$*

4. *Center and scale $x_i$, $y_i$*

We compute the coefficient of determination for all of the plots to assess the amount of linearity in each panel, computed as

$$R^2 = 1 - RSS/TSS, \qquad (1)$$

where TSS is the total sum of squares, $TSS = \sum_{i=1}^{N}(y_i - \bar{y})^2$ and $RSS = \sum_{i=1}^{N} e_i^2$, the residual sum of squares. The expected correlation coefficient $\rho$ in this scenario is

$$E[\rho] = \sqrt{\frac{\frac{1}{3}}{\frac{1}{3} + \sigma_T^2}},$$

because $E[RSS] = N\sigma_T^2$ and $E[TSS] = \sum_{i=1}^{N} E\left[y_i^2\right]$ (as $E[Y] = 0$), where

$$E\left[y_i^2\right] = E\left[x_i^2 + e_i^2 + 2x_i e_i\right] = \frac{1}{3} + \sigma_T^2.$$

The use of $R^2$ to assess the strength of the linear relationship (rather than the correlation, $r$) is indicated because human perception of correlation strength more closely aligns with $R^2$ (Bobko and Karren, 1979; Lewandowsky and Spence, 1989b).
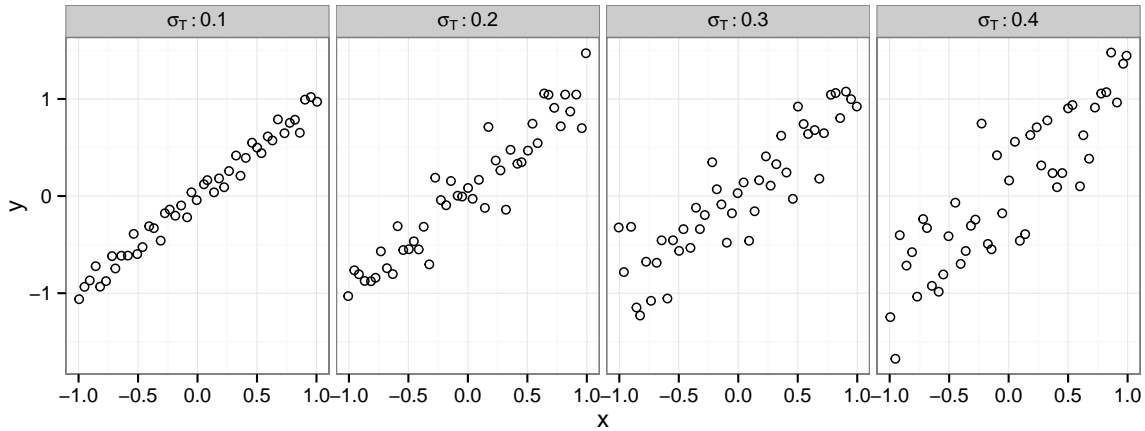


Figure 2: Set of scatterplots showing one draw each from the trend model $M_T$ for parameter values of $\sigma_T \in \{0.1, 0.2, 0.3, 0.4\}$.

### 2.1.2 Cluster Model $M_C$

We begin by generating $K$ cluster centers on a $K \times K$ grid, then we generate points around selected cluster centers.

**Algorithm 2.2**
*Input Parameters: N points, K clusters, $\sigma_C$ cluster standard deviation*
*Output: N points, in form of vectors x and y.*

1. *Generate cluster centers $(c_i^x, c_i^y)$ for each of the $K$ clusters, $i = 1, ..., K$:*

(a) in form of two vectors $c^x$ and $c^y$ of permutations of $\{1, ..., K\}$, such that

(b) the correlation between cluster centers $Cor(c^x, c^y)$ falls into a range of $[.25, .75]$.

2. Center and standardize cluster centers $(c^x, c^y)$:

$$\tilde{c}_i^x = \frac{c_i^x - \bar{c}}{s_c} \quad and \quad \tilde{c}_i^y = \frac{c_i^y - \bar{c}}{s_c},$$

where $\bar{c} = (K+1)/2$ and $s_c^2 = \frac{K(K+1)}{12}$ for all $i = 1, ..., K$.

3. For the $K$ clusters, we want to have nearly equal sized groups, but allow some variability. Group sizes are therefore determined as a draw from a multinomial distribution: determine group sizes $g = (g_1, ..., g_K)$, with $N = \sum_{i=1}^K g_i$, for clusters $1, ..., K$ as a random draw

$$g \sim Multinomial(K, p) \ where \ p = \tilde{p}/\sum_{i=1}^K \tilde{p}_i, \ for \ \tilde{p} \sim N(\frac{1}{K}, \frac{1}{2K^2}).$$

4. Generate points around cluster centers:

(a) $x_i = \tilde{c}_{g_i}^x + e_i^x$, where $e_i^x \sim N(0, \sigma_C^2)$

(b) $y_i = \tilde{c}_{g_i}^y + e_i^y$, where $e_i^y \sim N(0, \sigma_C^2)$

5. Center and scale $x_i$, $y_i$

As a measure of clustering we use a coefficient to assess the amount of variability within groups, compared to total variability. Note that for the purpose of clustering, variability is measured as the variability in both $x$ and $y$ from a common mean, i.e. we implicitly assume that the values in $x$ and $y$ are on the same scale (which we achieve by scaling in the final step of the generation algorithm).

> add cluster equation

(2)

> For the study we used $a = 1$, right?
> Susan: Yes, and I thought I'd purged all $a$ from the description accordingly, but I see that I missed one. It's fixed now.

### 2.1.3   Null Model $M_0$

The generative model for null data is a mixture model $M_0$ that draws $n_c \sim \text{Binomial}(N, \lambda)$ observations from the cluster model, and $n_T = N - n_c$ from the regression model $M_T$. Observations are assigned groups using hierarchical clustering, which creates groups consistent with any structure present in the generated data. This provides a plausible grouping for use in aesthetic and statistics requiring categorical data (color, shape, bounding ellipses).

Null data in this experiment is generated using $\lambda = 0.5$, that is, each point in a null data set is equally likely to have been generated from $M_C$ and $M_T$.
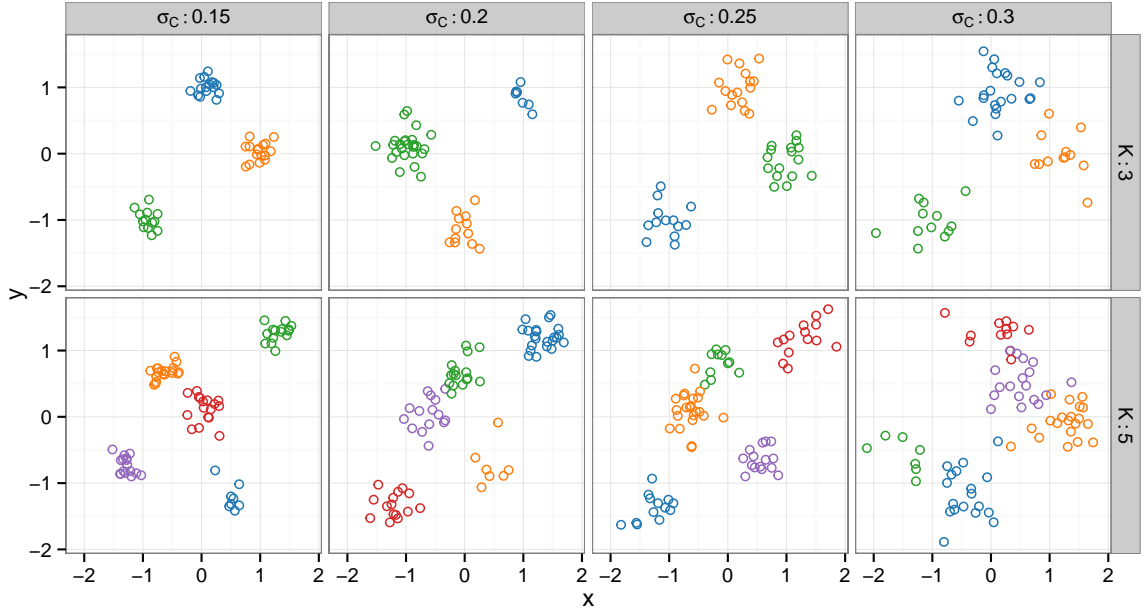
Figure 3: Scatterplots of clustering output for different inner cluster spread $\sigma_C$ (left to right) and different number of clusters $K$ (top and bottom).

### 2.1.4 Parameters used in Data Generation

These models provide the foundation for this experiment; by manipulating cluster standard deviation $\sigma_C$ and regression standard deviation $\sigma_T$ (directly related to correlation strength) for varying numbers of clusters $K = 3, 5$, we can systematically control the statistical signal present in the target plots and generate corresponding null plots that are mixtures of the two distributions. For each parameter set $\{K, N, \sigma_C, \sigma_T\}$, as described in table 1, we generate a lineup dataset consisting of one set drawn from $M_C$, one set drawn from $M_T$, and 18 sets drawn from $M_0$.

| Parameter | Description | Choices |
|---|---|---|
| $K$ | # Clusters | 3, 5 |
| $N$ | # Points | $15 \cdot K$ |
| $\sigma_T$ | Scatter around trend line | .15, .25, .35 |
| $\sigma_C$ | Scatter around cluster centers | .15, .20, .25 $(K = 3)$ <br> .20, .25, .30 $(K = 5)$ |

Table 1: Parameter settings for generation of lineup datasets.

> the simulation is so nice, maybe you could include one more sentence on the idea. Niladri did something similar in his Where's Waldo paper. Need to ask Di for a reference.
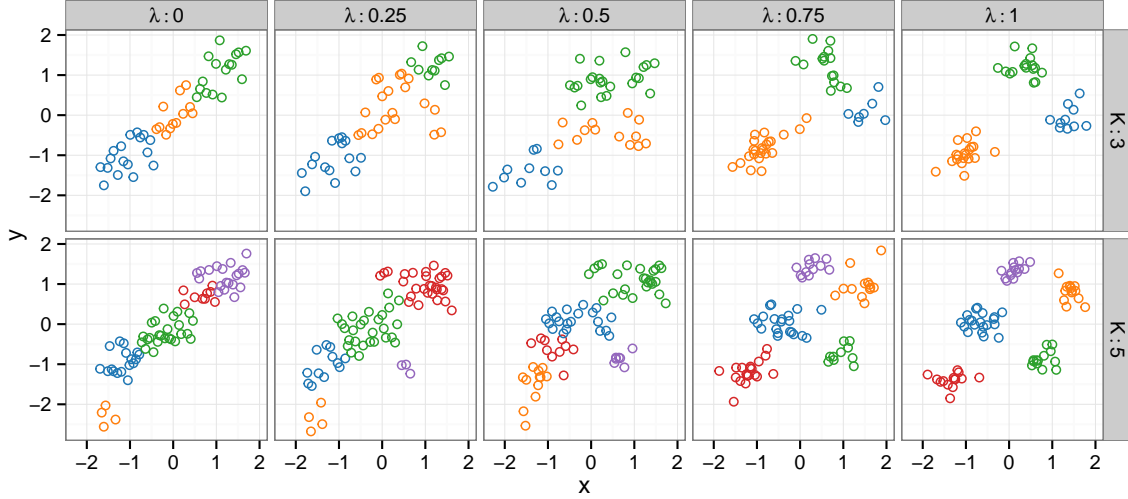
8

Figure 4: Scatterplots of data generated from $M_0$ using different values of $\lambda$.

The parameter values were chosen after examining the full parameter space through simulation of 1000 lineup datasets for each combination of $\sigma_T \in \{0.2, 0.25, ..., 0.5\}$, $\sigma_C \in \{0.1, 0.15, ..., 0.4\}$, and $K \in \{3, 5\}$; more complete results are provided in appendix 4. In accordance with the simulation, we identified values of $\sigma_T$ and $\sigma_C$ corresponding to "easy", "medium" and "hard" numerical comparisons between corresponding target data sets and null data sets. It is important to note that the numerical measures we have described in equations (1) and (2) only provide information on the numerical discriminability of the target datasets from the null datasets; the simulation cannot provide us with information on the perceptual discriminability, and it has been established that human perception of scatterplots does not replicate statistical measures exactly (Bobko and Karren, 1979; Mosteller et al., 1981; Lewandowsky and Spence, 1989b).

Each of the generated datasets is then plotted as a lineup, where we apply aesthetics which emphasize clusters and/or linear relationships, to experimentally determine how these aesthetics change participants' ability to identify each target plot. The next section describes the aesthetic combinations and their anticipated effect on participant responses.

## 2.2  Plot Aesthetics

Gestalt perceptual theory suggests that perceptual features such as shape, color, trend lines, and boundary regions modify the perception of ambiguous graphs, emphasizing clustering in the data (in the case of shape, color, and bounding ellipses) or linear relationships (in the case of trend lines and prediction intervals), as demonstrated in figure 1. For each dataset we examine the effect of plot aesthetics (color, shape) and statistical layers (trend line, boundary ellipses, prediction intervals) shown in table 2 on target identification. Examples of these plot aesthetics are shown in figure 5.
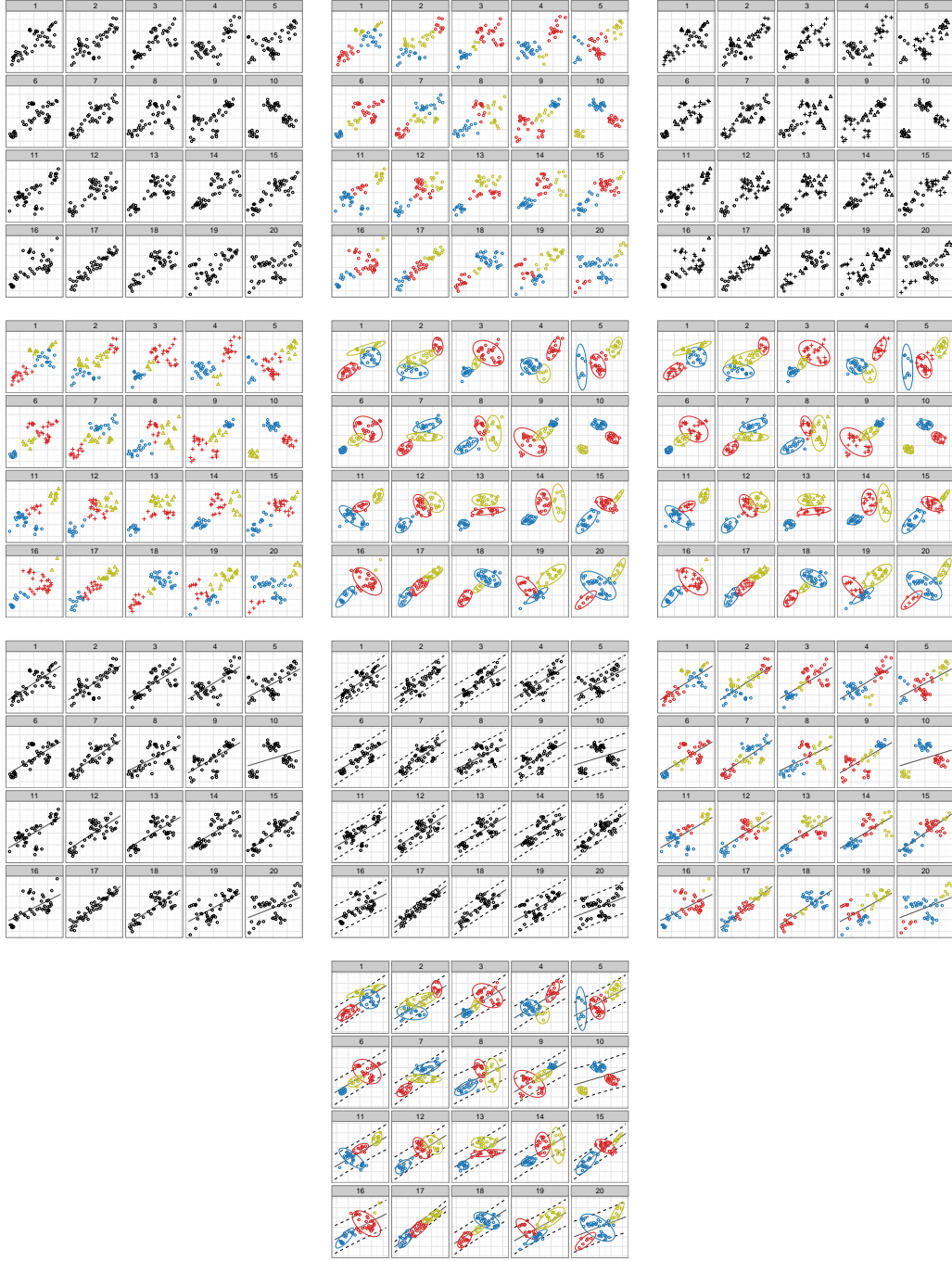
9

Figure 5: Each of the 10 plot feature combinations tested in this study, with $K = 3$, $\sigma_T = .25$ and $\sigma_C = .2$.

| | | Line Emphasis | | |
|---|---|---|---|---|
| | Strength | 0 | 1 | 2 |
| Cluster Emphasis | 0 | None | Line | Line + Prediction |
| | 1 | Color | Color + Line | |
| | | Shape | | |
| | 2 | Color + Shape | | Color + Ellipse + Line + Prediction |
| | | Color + Ellipse | | |
| | 3 | Color + Shape + Ellipse | | |

Table 2: Plot aesthetics and statistical layers which impact perception of statistical plots, according to gestalt theory.

We expect that relative to a plot with no extra aesthetics or statistical layers, the addition of color, shape, and 95% boundary ellipses increases the probability of a participant selecting the target plot with data generated from $M_C$, the cluster model, and that the addition of these aesthetics decreases the probability of a participant selecting the target plot with data generated from $M_T$, the linear model.

Similarly, we expect that relative to a plot with no extra aesthetics or statistical layers, the addition of a trend line and prediction interval increases the probability of a participant selecting the target plot with data generated from $M_T$, the linear model, and decreases the probability of a participant selecting the target plot with data generated from $M_C$, the cluster model.

## 2.3 Experimental Design

The study is designed hierarchically, as a factorial experiment for combinations of $\sigma_C$, $\sigma_T$, and $K$, with three replicates at each parameter combination. These parameters are used to generate lineup datasets which serve as blocks for the plot aesthetic level of the experiment; each dataset is rendered with every combination of aesthetics described in table 2. Participants are assigned to generated plots according to an augmented balanced incomplete block scheme: each participant is asked to evaluate 10 plots, which consist of one plot at each combination of $\sigma_C$ and $\sigma_T$, randomized across levels of $K$, with one additional plot providing replication of one level of $\sigma_C \times \sigma_T$. Each of a participant's 10 plots will present a different aesthetic combination.

Need to find some graphic/table which makes this a bit more clear.

## 2.4 Color and Shape Palettes

Colors and shapes used in this study were selected in order to maximize preattentive feature differentiation. Çağatay Demiralp et al. (2014) provide sets of 10 colors and 10 shapes, with corresponding distance matrices, determined by user studies. Using these perceptual kernels for shape and color, we identified sets of 3 and 5 colors and shapes which maximize the sum of pairwise differences, subject to certain constraints imposed by software and accessibility concerns.

The color palette used in Çağatay Demiralp et al. (2014) and shown in figure 6 is derived from colors available in Tableau visualization software.

Figure 6: Colors in Çağatay Demiralp et al. (2014). This study removed grey from the palette to make the experiment more inclusive of participants with colorblindness.
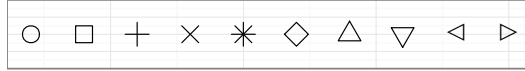


Figure 7: Shapes in Çağatay Demiralp et al. (2014). In order to control for varying point size due to Unicode vs. non-Unicode characters, the last two shapes were removed.

In order to produce experimental stimuli accessible to the approximately 4% of the population with red-green color deficiency (Gegenfurtner and Sharpe, 2001), we removed the grey hue from the palette. This modification produced maximally different color combinations which did not include red-green combinations, while also removing a color (grey) which is difficult to distinguish for those with color deficiency.

Software compatibility issues led us to exclude two shapes used in Çağatay Demiralp et al. (2014) and shown in figure 7. The left and right triangle shapes (available only in unicode within R) were excluded due to size differences between unicode and non-unicode shapes. After optimization over the sum of all pairwise distances, the maximally different shape sequences for the 3 and 5 group datasets also conform to the guidelines in Robinson (2003): for $K = 3$ the shapes are from Robinson's group 1, 2, and 9, for $K = 5$ the shapes are from groups 1, 2, 3, 9, and 10. Robinson's groups are designed so that shapes in different groups show differences in preattentive properties; that is, they are easily distinguishable. In addition, all shapes are non-filled shapes, which means that they are consistent with one of the simplest solutions to overplotting of points in the tradition of Tukey (1977); Cleveland (1994) and Few (2009). For this reason we abstained from the additional use of alpha-blending of points to diminish the effect of overplotting in the plots.

## 2.5   Hypotheses

The primary purpose of this study is to understand how visual aesthetics affect signal detection in the presence of competing signals. We expect that plot modifications which emphasize similarity and proximity, such as color, shape, and 95% bounding ellipses, will increase the probability of detecting the clustering relationship, while plot modifications which emphasize good continuation, such as trend lines and prediction intervals, will increase the probability of detecting the linear relationship.

A secondary purpose of the study is to relate signal strength (as determined by dataset parameters $\sigma_C$, $\sigma_T$, and $K$) to signal detection in a visualization by a human observer.
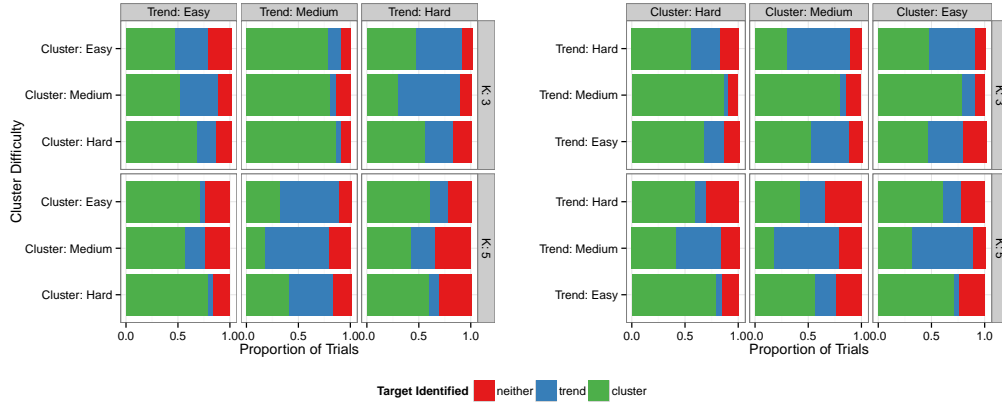
Figure 8: Identified targets for each level of parameter settings (across plot types).

## 2.6 Participant Recruitment

describe amazon turk, participant instructions, screening procedures, etc.

# 3 Results

Data collection was conducted over a 25 hour period, during which time 1380 individuals completed 13502 unique lineup evaluations (the average individual evaluated 9.78 lineups). Each plot was evaluated by between 12 and 40 individuals (Mean: 25.00, SD= 4.94). 82.2% of the participant evaluations identified at least one of the two target plots successfully (Trend: 26.5%, Cluster: 56.3%). Only 3% of participant evaluations identified more than one target plot, and of these multiple identifications, 21.5% identified both targets correctly.

I don't know where the demographic information for the experiment is stored...

## 3.1 Graphical Exploration

Figure 8 shows aggregate parameter-level accuracy rates. Participants appear to be less accurate when evaluating plots with $K = 3$ clusters, and appear to be at least as successful at evaluating "hard" cluster lineups as "easy" cluster lineups, conditional on trend difficulty. There also appears to be an interaction between number of clusters, $K$, and trend difficulty that may merit further evaluation.

Figure 9 shows aggregate accuracy rates for each plot aesthetic combination. It is again apparent that the cluster targets were overall more likely to be identified than line targets across all aesthetic combinations, however, it is also evident that plot aesthetics influence the identified target.
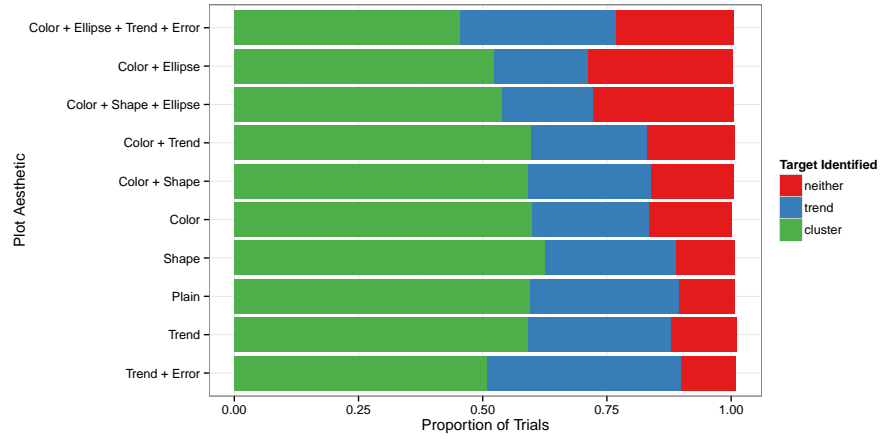
13

Figure 9: Proportion of trials identifying each target for each plot type (across parameter settings).
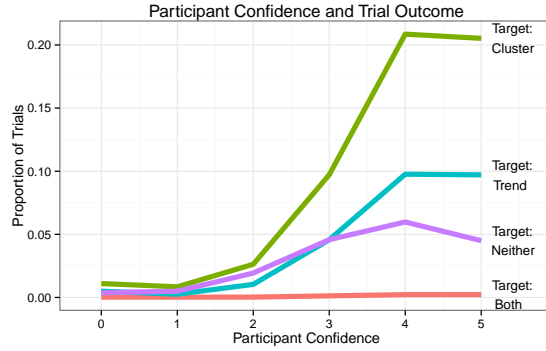


Figure 10: Participant confidence levels compared with trial results.

In addition to participant identification of target plots, we also asked that they rate their confidence in their answer. Figure 10 shows aggregate participant confidence rating as a function of trial outcome. Participants who did not identify either target plot were less likely to be "extremely confident" in their answer, while participants who identified either the trend or the cluster target correctly were highly confident that their answer was correct. Overall, though, participants seem to have some degree of confidence in their answer, regardless of whether the answer was correct.
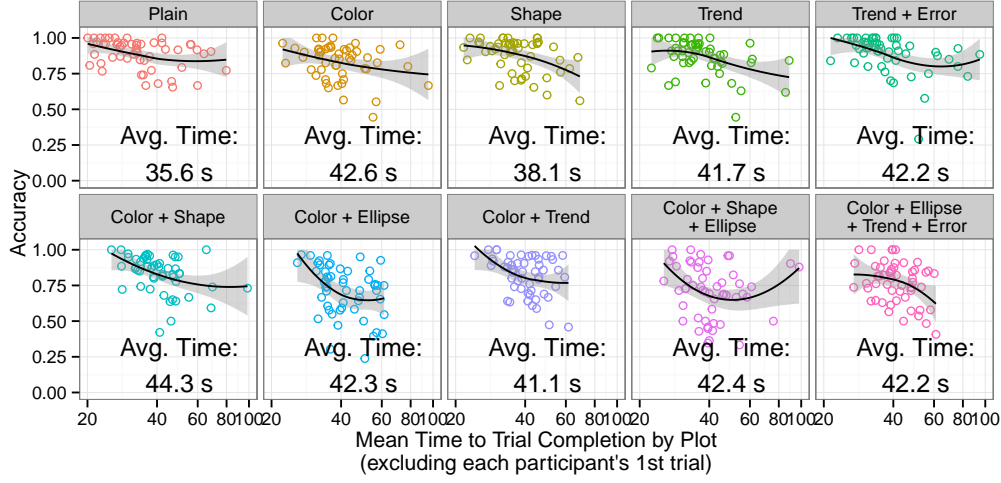
Figure 11: Accuracy (identifying either target plot) compared with mean trial time, by plot aesthetic.

As data collection was conducted entirely online, we cannot measure responses in the millisecond range characteristic of many psychometric studies, however, the data server does record the trial start and end time. Examining differences in average response times across trials provides us with an additional measure of trial difficulty or perceptual complexity. We can also explore whether participants spent more time on certain types of plots and whether that additional time increased accurate target identification. In order to remove the "novelty" effect of an unfamiliar task, we excluded every participant's first trial from this portion of the analysis. Average trial times and accuracy rates for each lineup are shown in figure 11 facetted by plot type, and in figure 12 facetted by parameter difficulty.

We will first consider the effect of plot aesthetics on target selection for each target type (separately), and then we will analyze the effect of parameter values on participant performance.

## 3.2 Linear Target Model

We will model the probability of selecting the linear target plot as a function of plot type, with random effects for dataset (which encompasses parameter effects) and participant (accounting for variation in individual skill level). For plot type $i = 1, ..., 10$ displaying dataset $j = 1, ..., 54$ by participant $k = 1, ..., P$,
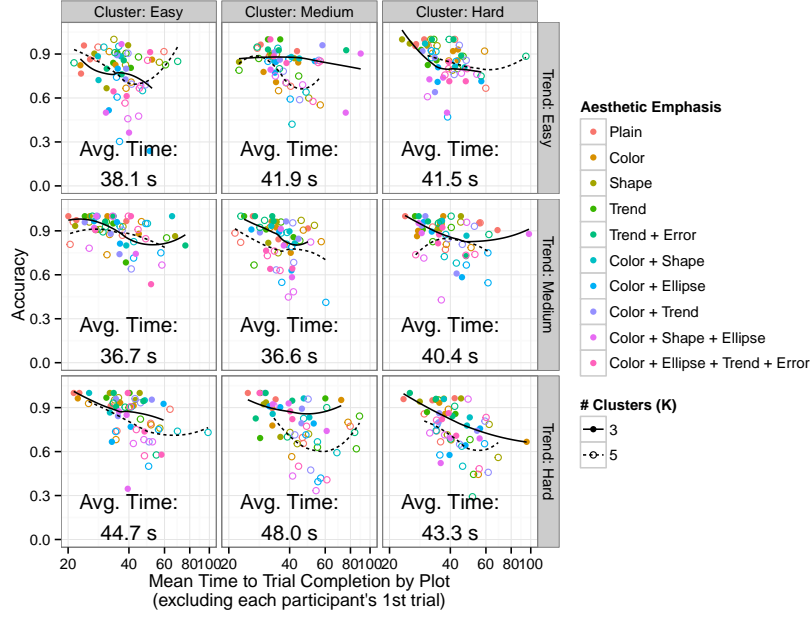
Figure 12: Accuracy (identifying either target plot) compared with mean trial time, by parameter settings.

$$P(\text{success}) = \left(e^{\theta}\right) / \left(1 + e^{\theta}\right) \tag{3}$$
$$\theta = \mathbf{X}\beta + \mathbf{J}\gamma + \mathbf{K}\eta + \epsilon \tag{a}$$

where $\beta_i$ describe the effect of specific plot aesthetics

$$\gamma_j \overset{iid}{\sim} N\left(0, \sigma_{\text{data}}^2\right), \text{ the random effect for dataset specific characteristics}$$

$$\eta_k \overset{iid}{\sim} N\left(0, \sigma_{\text{participant}}^2\right), \text{ the random effect for participant characteristics}$$

and $\epsilon_{ijk} \overset{iid}{\sim} N\left(0, \sigma_e^2\right)$, the error associated with a single trial evaluation

We note that any variance due to parameters $K$, $\sigma_T$, and $\sigma_C$ is contained within $\sigma_{\text{data}}^2$ and can be examined using a subsequent model.
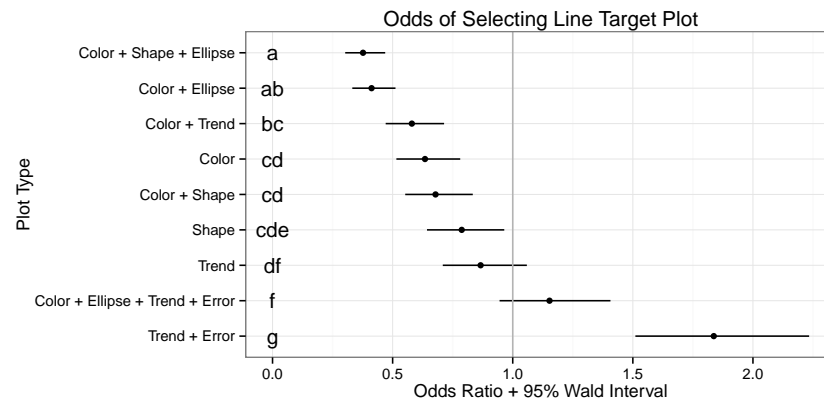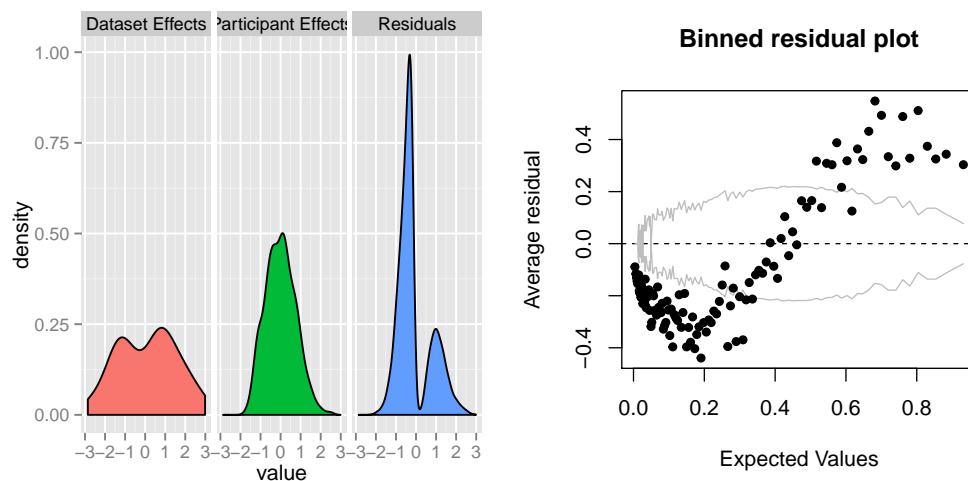
Figure 13: Odds ratios describing the odds of detecting the linear target plot for each aesthetic, relative to a plain scatterplot. Only the combination of Trend + Error significantly increases the odds of linear target plot detection relative to the control plot.



I'm a bit concerned about the residuals here... thoughts? I'm also not sure where else to go with this particular discussion... yes, the effects exist. No, the dataset effects don't look all that normal, because there's a bunch of parameters in there. Participant effects are normal-ish, though, so that's something.

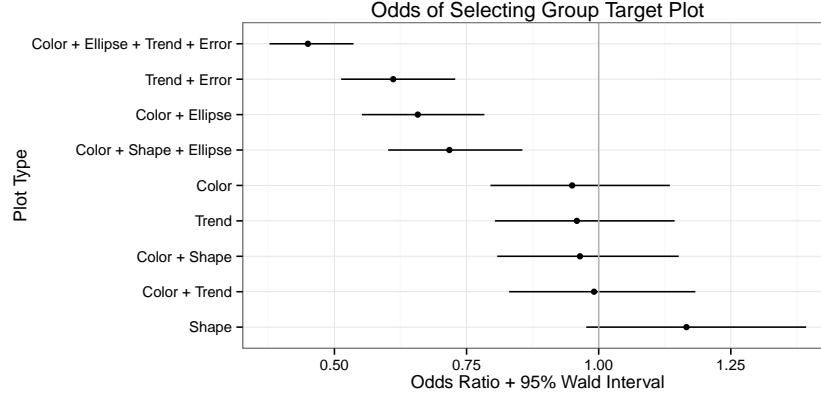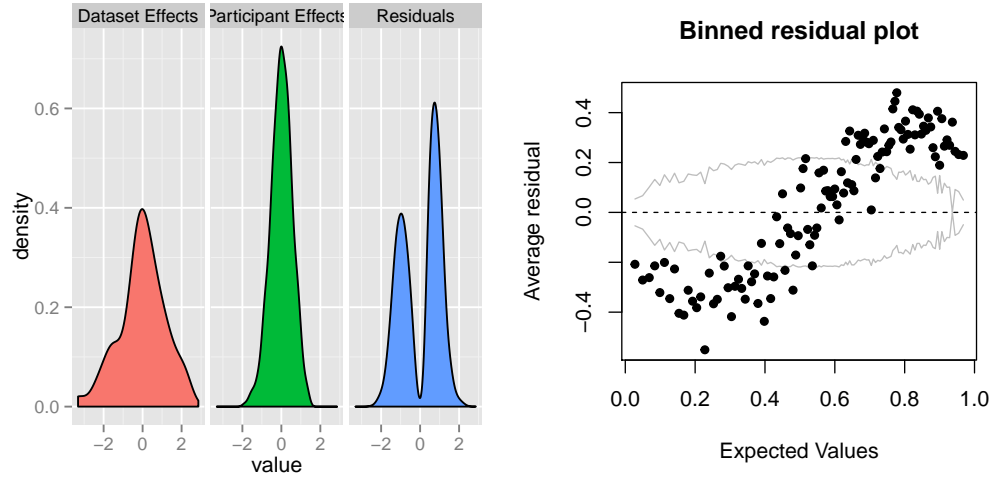Discuss variance/covariance and random effects??

Figure 14: Odds ratios describing the odds of detecting the cluster target plot for each aesthetic, relative to a plain scatterplot. The presence of error lines or bounding ellipses significantly decreases the probability of correct target detection, and no aesthetic successfully increases the probability of correct target detection. This may be due to differences in group size for null plots, with data generated under $M_0$ compared with the group target plot displaying data generated under $M_C$.

## 3.3 Group Target Selection

We now examine the probability of selecting the group target plot as a function of plot type, with random effects for dataset (which encompasses parameter effects) and participant (accounting for variation in individual skill level). The model fit here is the same as that shown in equation (3), except that success in this model is defined as identification of the cluster target plot.



18

## 3.4 Face-Off: Group versus Line

Just another idea of evaluating this data set: For each data set we only consider those evaluations that correctly identify one of the targets. This reveals that participants overall favored groups to lines at a ratio of about 2:1. We remove this overall effect using an intercept, and model group vs line decisions using a logistic regression with a random effect for each dataset to account for different difficulty levels in the generated data. The estimated odds of a decision in favor of group over line target are shown in figure 15. From left to right the odds of selecting the group target over the line target increase. As hypothesized, the strongest signal for identifying groups, is color + shape + ellipse, while trend + error results in the strongest signal in favor of trends. Most of the effects are not significantly different (see the letter values Piepho (2004) on the left hand side of the figure, representing pairwise comparisons of all of the designs, adjusted for multiple comparison). Trend+error plots and ellipse+trend plots are significantly different from all of the other designs. Apart from that, the only significant difference between designs is between color+shape+ellipse plots and trend plots. Surprisingly(?), the two designs sending mixed signals (color + trend, color+ ellipse+trend) end on opposite ends of the scale.

My take on this:
The similarity/proximity effect (as indicated by clustering and color/shape/etc.) dominates the equation, including dominating the color+trend (good continuation) condition.
The addition of common region (ellipses, error bars) modifies this effect somewhat, reinforcing the clustering and trend when present (as those are at the extreme ends of the plot).
When (trend + error) are present in the same plot, you get additional gestalt ordering principle(s): common region + good continuation **+ connectedness + closure**:
since the error band would be perceived as containing the points and connected to the trendline - this does somewhat lessen with the lines for error bands we have shown here, but I think the point still holds.
Plus, you could argue for closure, since you have fairly symmetric lines on either side of a central object.
(color + ellipse) + clustering = (similarity + common region) + proximity
is not as strong as
(trend + error) + correlation = (good continuation (trendline) + common region + connectedness + closure) + good continuation (points)

## 3.5 Signal Strength

# 4 Discussion

# References

Bobko, P. and Karren, R. (1979), "The perception of Pearson product moment correlations from bivariate scatterplots," *Personnel Psychology*, 32, 313–325.

Çağatay Demiralp, Bernstein, M., and Heer, J. (2014), "Learning Perceptual Kernels for Visualization Design," *IEEE Trans. Visualization & Comp. Graphics (Proc. InfoVis)*.
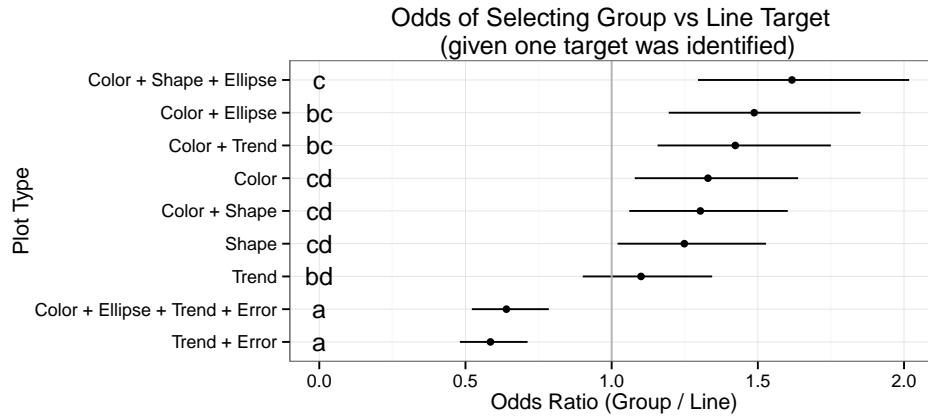
Figure 15: Estimated odds of decision for group versus line target based on evaluations that resulted in the identification of one of these targets. Plot types are significantly different, if they do not share a letter as given on the left hand side of the plot.

Cleveland, W. S. (1994), *The Elements of Graphing Data*, Hobart Press, 1st ed.

Cleveland, W. S. and McGill, R. (1984), "Graphical Perception: Theory, Experimentation, and Application to the Development of Graphical Methods," *Journal of the American Statistical Association*, 79, pp. 531–554.

Few, S. (2009), *Now You See It: Simple Visualization Techniques for Quantitative Analysis*, Burlingame, CA: Analytics Press, 1st ed.

Gegenfurtner, K. R. and Sharpe, L. T. (2001), *Color vision: From genes to perception*, Cambridge University Press.

Goldstein, E. B. (2009), *Encyclopedia of perception*, Sage Publications.

Healey, C. G. and Enns, J. T. (1999), "Large datasets at a glance: Combining textures and colors in scientific visualization," *Visualization and Computer Graphics, IEEE Transactions on*, 5, 145–167.

— (2012), "Attention and visual memory in visualization and computer graphics," *Visualization and Computer Graphics, IEEE Transactions on*, 18, 1170–1188.

Lewandowsky, S. and Spence, I. (1989a), "Discriminating strata in scatterplots," *Journal of the American Statistical Association*, 84, 682–688.

— (1989b), "The perception of statistical graphs," *Sociological Methods & Research*, 18, 200–242.

Majumder, M., Hofmann, H., and Cook, D. (2013), "Validation of Visual Statistical Inference, Applied to Linear Models," *Journal of the American Statistical Association*, 108, 942–956.

Mosteller, F., Siegel, A. F., Trapido, E., and Youtz, C. (1981), "Eye fitting straight lines," *The American Statistician*, 35, 150–152.

Piepho, H.-P. (2004), "An algorithm for a letter-based representation of all-pairwise comparisons," *Journal of Computational and Graphical Statistics*, 13, 456–466.

Robinson, H. (2003), "Usability of Scatter Plot Symbols," *ASA Statistical Computing & Graphics Newsletter*, 14, 9–14.

Spence, I. and Garrison, R. F. (1993), "A remarkable scatterplot," *The American Statistician*, 47, 12–19.

Treisman, A. (1985), "Preattentive processing in vision," *Computer Vision, Graphics, and Image Processing*, 31, 156 – 177.

Tukey, J. W. (1977), *Exploratory Data Analysis*, Lebanon, IN: Addison Wesley.

# Simulation Studies of Parameter Space

## 4.1  Distribution of Test Statistics

Simulating lineup data sets, we can compare test statistics measuring trend strength, cluster strength, and cluster size inequality for the null plots and target plots. These distributions allow us to objectively assess the difficulty of detecting the target datasets computationally (without relying on human perception).

Add equations for test statistics

Figure 16 show computed densities of the maximum null distribution measure compared with the measure in the signal plot. There is some overlap in the distribution of $R^2$ for the null plots compared with the target plot displaying data drawn from $M_T$. We have two measures comparing data drawn from $M_C$ and $M_0$; the cluster measure examines the variance in $x$ and $y$ described by the cluster center; the gini coefficient examines the inequality in group sizes. These simulations indicate that it may be possible to differentiate $M_C$ based on two different features in clustered data. In future experiments, it may be beneficial to control cluster size more tightly to remove this additional feature.

The distribution of the cluster statistic values are more easily separated from the null plots than the distribution of the line statistic, indicating that $\sigma_C = 0.20$ is producing target plots that are a bit easier to spot than trend targets with a parameter value of $\sigma_T = 0.25$, however, the inequality of group sizes may distract participants from the intended target signal of cluster cohesion.

## 4.2  Full Parameter Space Simulation Study

Using 1000 simulations for each of the 98 combinations of parameters ($K = \{3, 5\}$, $\sigma_C = \{.1, .15, .2, .25, .3, .35, .4\}$, $\sigma_T = \{.2, .25, .3, .35, .4, .45, .5\}$), we explored the effect of parameter value on the distribution of summary statistics describing the line strength ($R^2$) and cluster strength for null and target plots.

Describe statistics

Figures 17 and 18 show the 25th and 75th percentiles of the distribution of line and cluster summary statistics for each set of parameter values. These plots guided our evaluation of "easy", "medium" and "hard" parameter values for line and cluster tasks.
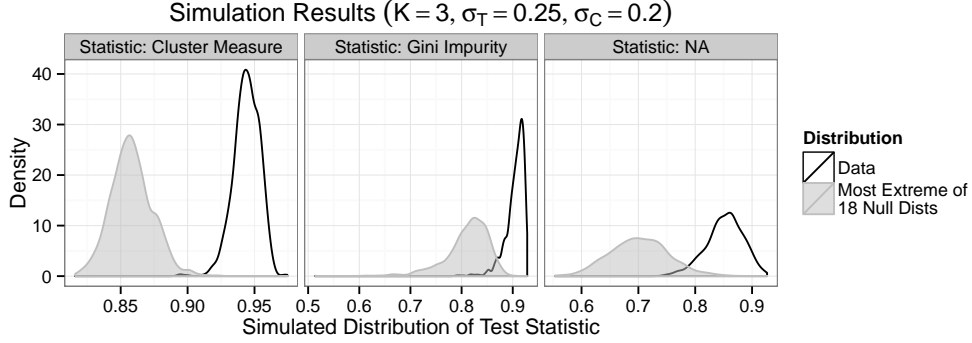
Figure 16: Density of test statistics measuring trend strength, cluster strength, and cluster inequality for target distributions and null plots.

Additionally, we note that there is an interaction between $\sigma_C$ and $\sigma_T$: the distinction between target and null on a fixed setting of clustering becomes increasingly difficult as the standard deviation for the linear trend is increased, and vice versa. There may additionally be a three-way interaction between $\sigma_C, \sigma_T$, and $K$: the size of the blue intervals (bottom figure) changes in size between different levels of $K$, it changes for different levels of $\sigma_C$ and $\sigma_T$. These interactions suggest that in order to examine differences in aesthetics, we must block by parameter settings (this can be accomplished through blocking by dataset). Each dataset is non-deterministic, because we have a random process generating from different parameter settings, not a deterministic run setting as in an engineering setting. It is thus important to use replicates of each parameter setting to ensure that we can separate data-level effects from parameter-level effects.

Additionally, after the experiment was complete, we examined the distribution of group size (as measured by gini impurity) to establish whether there were any systematic differences in group size inequality between data generated from $M_0$ (null data) and data generated from $M_C$ (cluster data). Figure 19 demonstrates that the cluster plots have significantly lower group size differences than null plots at all parameter combinations. It is therefore possible that some participants will identify extraordinarily unequal group sizes present in null plots as significantly different from the other lineup plots, ignoring any cluster signal. Future studies should more tightly control group size in order to reduce this effect.

## Experimental Design

Initially, assume a fully factorial, balanced design, with $r$ unique datasets per parameter set (replicates) and $P$ evaluations per (aesthetic|dataset). The experiment is conducted at three levels: parameter sets (with replication, so EUs are data sets), plot types (i.e. a certain set of aesthetics), and participant evaluations. At the first level, there are three parameters: $K = 3$, $\sigma_T \in \{.15, .25, .35\}$, and $\sigma_C \in \{.15, .25, .35\}$. At the second level, there are blocks (by data set), and then 4 aesthetic combinations.
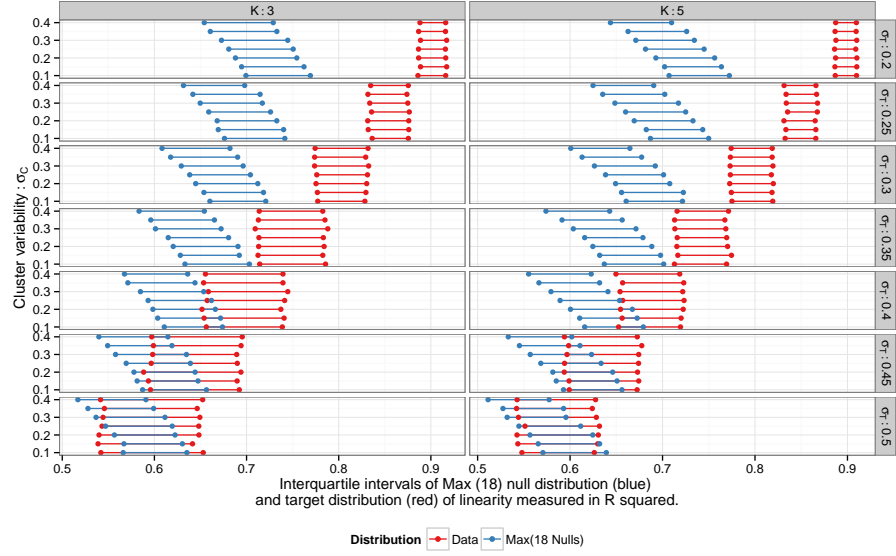
Figure 17: Simulated interquartile range of $R^2$ values for target and null data distributions.
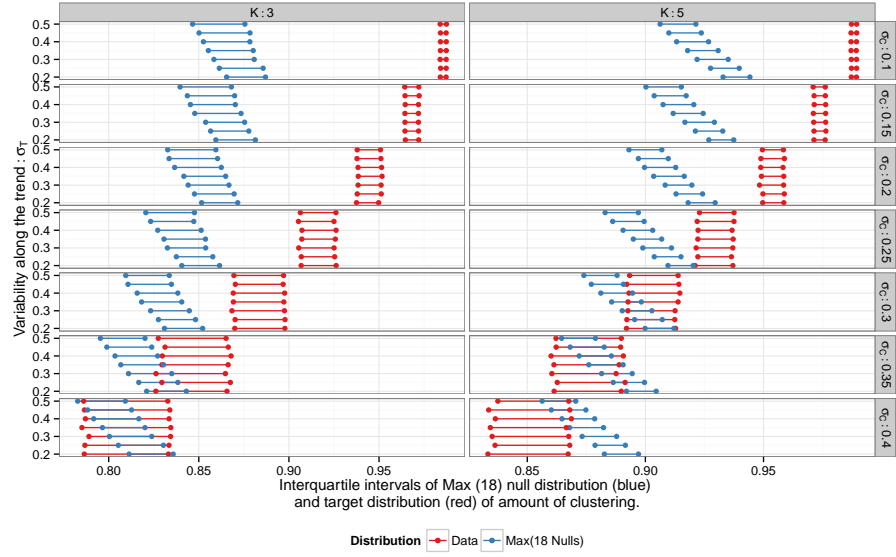


Figure 18: Simulated interquartile range of cluster cohesion statistic values for target and null data distributions.
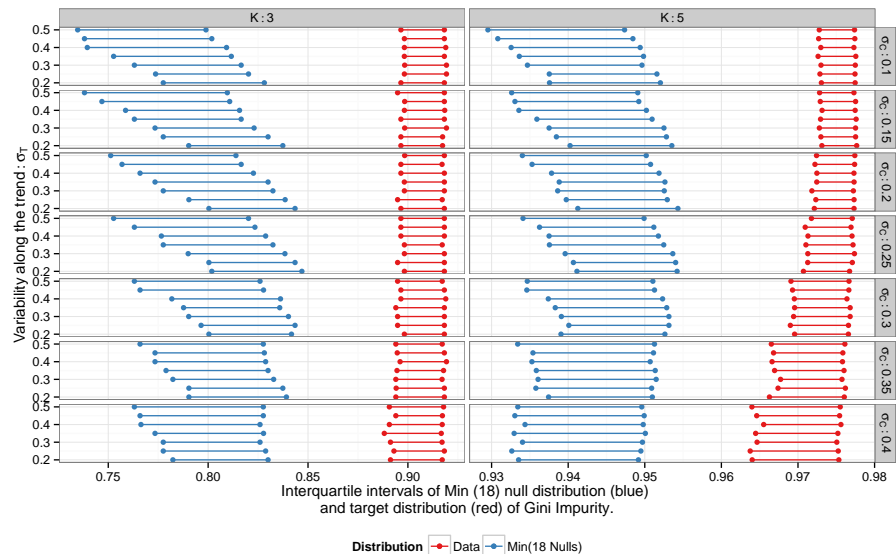
Figure 19: Simulated interquartile range of group size inequality statistic values for cluster and null data distributions.

We'll have to use contrasts to measure the effect of color individually, etc., for now let's just consider the ANOVA evaluation

Finally, at the lowest level, there are participant effects.

At the participant level, we need to decide if we're going to fully randomize, try to block, etc. - are participants going to get 10 different data sets? 5? Not sure how to conceptualize that, and I would imagine it will affect how we organize model evaluation. Grr, I hate mixed models. HH: yes, I would assume that participants get ten plots each, one from each of the designs in a random order. (We have the data base set up that way).

Modified from Table 10.6 (pg 181) of Design of Experiments by Dr. Morris. The table in the book has a four-factor split plot design with three levels (randomized, block, block).

We have a couple of options:

- keep the full factorial experiment, use one (at most two) replicates, and use higher level factorial effects to beef up any error variance terms.

- Do a full factorial experiment for $K = 3$ and use a subset of the factorial experiment for $K = 5$ (either using a subset of cases for $\sigma_T$ and $\sigma_C$, or a subset of combinations of the two cases/fractional factorial.)

| Level | Factor | Source | DF | Sum of Squares |
|---|---|---|---|---|
| | $K$ | $\alpha$ | 1 | $\sum_i(4)(5)(r)(10P)(\overline{y}_{i.....} - \overline{y}_{......})^2$ |
| | $\sigma_T^2$ | $\beta$ | 3 | $\sum_j(2)(5)(r)(10P)(\overline{y}_{.j....} - \overline{y}_{......})^2$ |
| | $\sigma_C^2$ | $\gamma$ | 4 | $\sum_i(2)(4)(r)(10P)(\overline{y}_{..k...} - \overline{y}_{......})^2$ |
| Dataset | | $(\alpha\beta)$ | 3 | $\sum_{ij}(5)(r)(10P)(\overline{y}_{ij....} - \overline{y}_{i.....} - \overline{y}_{.j....} + \overline{y}_{......})^2$ |
| | | $(\alpha\gamma)$ | 4 | $\sum_{ik}(4)(r)(10P)(\overline{y}_{i.k...} - \overline{y}_{i.....} - \overline{y}_{..k...} + \overline{y}_{......})^2$ |
| | | $(\beta\gamma)$ | 12 | $\sum_{jk}(2)(r)(10P)(\overline{y}_{.jk...} - \overline{y}_{.j....} - \overline{y}_{..k...} + \overline{y}_{......})^2$ |
| | | $(\alpha\beta\gamma)$ | 12 | $\sum_{ijk}(r)(10P)(\overline{y}_{ijk...} - \overline{y}_{i.....} - \overline{y}_{.j....} - \overline{y}_{..k...}$ $+\overline{y}_{ij....} + \overline{y}_{.jk...} + \overline{y}_{i.k...} - \overline{y}_{......})^2$ |
| | | Resid. | $(2)(4)(5)(r-1)$ | $\sum_{ijkl}(10P)(\overline{y}_{ijkl..} - \overline{y}_{i.....} - \overline{y}_{.j....} - \overline{y}_{..k...} + \overline{y}_{ij....} + \overline{y}_{.jk...}$ $+\overline{y}_{i.k...} - \overline{y}_{ijk...} - \overline{y}_{ij.l..} - \overline{y}_{i.kl..} - \overline{y}_{.jkl..} + \overline{y}_{......})^2$ |
| | Total | | $(2)(4)(5)(r)-1$ | $\sum_{ijkl}(10P)(\overline{y}_{ijkl..} - \overline{y}_{......})^2$ |
| | Dataset | blocks | $(2)(4)(5)(r)-1$ | $\sum_{ijkl}(10P)(\overline{y}_{ijkl..} - \overline{y}_{......})^2$ |
| | Aes. | $\delta$ | 9 | $\sum_m(2)(4)(5)(P)(\overline{y}_{....m.} - \overline{y}_{......})^2$ |
| | Aes x $K$ | $(\alpha\delta)$ | 9 | $\sum_{im}(4)(5)(P)(\overline{y}_{i...m.} - \overline{y}_{i.....} - \overline{y}_{....m.} + \overline{y}_{......})^2$ |
| Plot | Aes x $\sigma_T$ | $(\beta\delta)$ | 27 | $\sum_{jm}(2)(5)(P)(\overline{y}_{.j..m.} - \overline{y}_{.j....} - \overline{y}_{....m.} + \overline{y}_{......})^2$ |
| | Aes x $\sigma_C$ | $(\gamma\delta)$ | 36 | $\sum_{km}(2)(4)(P)(\overline{y}_{..k.m.} - \overline{y}_{..k...} - \overline{y}_{....m.} + \overline{y}_{......})^2$ |
| | Others | | 9(31) | difference |
| | Resid | | 40(rP-1)-(40r-1) | |
| | Total | | $400r-1$ | $\sum_{ijklm}(P)(\overline{y}_{ijklm.} - \overline{y}_{ijkl..})^2$ |
| | Picture | Sub-blocks | $400r-1$ | $\sum_{ijklm}(P)(\overline{y}_{ijklm.} - \overline{y}_{ijkl..})^2$ |
| Trial | Participants | $\tau$ | $P-1$ | $\sum_n(2)(4)(5)(r)(10)(\overline{y}_{.....n} - \overline{y}_{......})^2$ |
| | Resid | | $(400r-1)(P)$ | difference |
| | Total | | $400(r)(P)-1$ | $\sum_{ijklmn}(y_{ijklmn} - \overline{y}_{......})^2$ |

Table 3: Evaluation of sources of error in a full factorial version of the experiment, with $r$ replicates of each parameter combination and $P$ participant evaluations of each plot(data/aesthetic combination).

The fractional factorial option will be a pain to explain when we write things up; it will be simpler to explain using a subset of cases. Given that we don't particularly care about the third-order effects (and possibly not even the second-order effects) for the parameters, I'm inclined to say that the single-replicate option is the easiest way to go (and lets us keep the simple SSQ in the table, which is a huge bonus in my opinion). Even if we just use the third-order interaction effect as error, we still have 12 degrees of freedom; that should be plenty - we'd only need F=2.69 to get a significant result for even the $(\sigma_T\sigma_C)$ test.

HH: We might not care about interpreting the two-way interactions, but unfortunately they will be there (see comment at the back). So I would suggest to go with a full factorial design in $\sigma_C, \sigma_T$, and $K$, with three replications each (we need the replicates, also explained in the back). This gives us 18 parameter settings, and $18 \cdot 3 = 54$ data sets. In case you still want to consider the effect of the number of datapoints $N$, we could switch from fully factorial to fractional factorial and replace the three-way interaction of $\sigma_C, \sigma_T$, and $K$ by the settings of $N$. That way we will keep the 18 settings.

# Model Results

Table 4: ANOVA table - only one replicate. Evaluation of sources of error in a full factorial version of the experiment, with one replicate of each parameter combination and $P$ participant evaluations of each plot(data/aesthetic combination).

| Level | Factor | Source | DF | Sum of Squares |
|---|---|---|---|---|
| Dataset | $K$ | $\alpha$ | 1 | $\sum_i (4)(5)(10P)(\overline{y}_{i....} - \overline{y}_{.....})^2$ |
| | $\sigma_T^2$ | $\beta$ | 3 | $\sum_j (2)(5)(10P)(\overline{y}_{.j...} - \overline{y}_{.....})^2$ |
| | $\sigma_C^2$ | $\gamma$ | 4 | $\sum_i (2)(4)(10P)(\overline{y}_{..k..} - \overline{y}_{.....})^2$ |
| | | Resid. | 22 | difference |
| | Total | | 39 | $\sum_{ijk}(10P)(\overline{y}_{ijk..} - \overline{y}_{.....})^2$ |
| Plot | Dataset | blocks | 39 | $\sum_{ijk}(10P)(\overline{y}_{ijk..} - \overline{y}_{.....})^2$ |
| | Aes. | $\delta$ | 9 | $\sum_m (2)(4)(5)(P)(\overline{y}_{...m.} - \overline{y}_{.....})^2$ |
| | Aes x $K$ | $(\alpha\delta)$ | 9 | $\sum_{im}(4)(5)(P)(\overline{y}_{i..m.} - \overline{y}_{i....} - \overline{y}_{...m.} + \overline{y}_{.....})^2$ |
| | Aes x $\sigma_T$ | $(\beta\delta)$ | 27 | $\sum_{jm}(2)(5)(P)(\overline{y}_{.j.m.} - \overline{y}_{.j...} - \overline{y}_{...m.} + \overline{y}_{.....})^2$ |
| | Aes x $\sigma_C$ | $(\gamma\delta)$ | 36 | $\sum_{km}(2)(4)(P)(\overline{y}_{..km.} - \overline{y}_{..k..} - \overline{y}_{...m.} + \overline{y}_{.....})^2$ |
| | Resid | | 9(31) | difference |
| | Total | | 399 | $\sum_{ijkm}(P)(\overline{y}_{ijkm.} - \overline{y}_{ijk..})^2$ |
| Trial | Picture | Sub-blocks | 399 | $\sum_{ijkm}(P)(\overline{y}_{ijkm.} - \overline{y}_{ijk..})^2$ |
| | Participants | $\tau$ | $P-1$ | $\sum_n (2)(4)(5)(10)(\overline{y}_{....n} - \overline{y}_{.....})^2$ |
| | Resid | | $399(P-1)$ | difference |
| | Total | | $400P-1$ | $\sum_{ijkmn}(y_{ijkmn} - \overline{y}_{....})^2$ |

| Plot Aesthetic | Log Odds | Std. Error | Z | P value | Tukey Post Hoc Differences |
|---|---|---|---|---|---|
| Trend + Error | 0.6081 | 0.0998 | 6.09 | 0.0000 | g |
| Color + Ellipse + Trend + Error | 0.1425 | 0.1015 | 1.40 | 0.1603 | f |
| Trend | -0.1436 | 0.1024 | -1.40 | 0.1608 | df |
| Shape | -0.2387 | 0.1035 | -2.31 | 0.0211 | cde |
| Color + Shape | -0.3877 | 0.1050 | -3.69 | 0.0002 | cd |
| Color | -0.4546 | 0.1061 | -4.29 | 0.0000 | cd |
| Color + Trend | -0.5448 | 0.1062 | -5.13 | 0.0000 | bc |
| Color + Ellipse | -0.8861 | 0.1104 | -8.03 | 0.0000 | ab |
| Color + Shape + Ellipse | -0.9763 | 0.1118 | -8.73 | 0.0000 | a |

Table 5: Fitted values of fixed effects for the model described in (3). Only Trend+Error plots significantly increase the probability of detecting the linear target plot (with data generated from $M_T$), while most other aesthetic combinations decrease the probability of detecting the linear target plot.

| Plot Aesthetic | Log Odds | Std. Error | Z | P value | Tukey Post Hoc Differences |
|---|---|---|---|---|---|
| Shape | 0.1536 | 0.0906 | 1.70 | 0.0898 | d |
| Color + Trend | -0.0088 | 0.0902 | -0.10 | 0.9220 | d |
| Color + Shape | -0.0361 | 0.0903 | -0.40 | 0.6897 | d |
| Trend | -0.0422 | 0.0900 | -0.47 | 0.6395 | d |
| Color | -0.0514 | 0.0907 | -0.57 | 0.5705 | cd |
| Color + Shape + Ellipse | -0.3320 | 0.0899 | -3.69 | 0.0002 | bc |
| Color + Ellipse | -0.4191 | 0.0895 | -4.68 | 0.0000 | b |
| Trend + Error | -0.4924 | 0.0896 | -5.49 | 0.0000 | b |
| Color + Ellipse + Trend + Error | -0.7992 | 0.0898 | -8.90 | 0.0000 | a |

Table 6: Fitted values of fixed effects for the model described in (**??**).