

Beautification of Images by Generative Adversarial Networks (GANs)

Amar Music

Master of Psychology: Theory and Research

KU LEUVEN

Faculty of Psychology and Educational Sciences

Supervisor: Prof. Johan Wagemens

Co-supervisor: Anne-Sofie Maerten

July 20, 2022

Contents

1	Abstract	2
2	Beautification of images by Generative Adversarial Networks (GANs)	3
2.1	Defining Beauty	3
2.2	Machine Learning and Beauty	4
3	Methods	5
3.1	Participants	5
3.2	Questionnaire	5
3.3	Behavioral Experiment	5
3.4	Data Analysis	8
3.4.1	Structural Equation Modeling	8
3.4.2	Behavioral Task	9
4	Results	10
4.1	Structural Models	10
4.2	Behavioral Task	12
4.2.1	Descriptive	12
4.2.2	Psychometric Function	12
4.2.3	Confounding Factors	12
5	Discussion	13
5.1	Behavioral Experiment	13
5.2	GANalyze Parameters	13
6	Appendix	16
6.1	Appendix A: Shortened Aesthetic Experience Questionnaire	16
6.2	Appendix B: Task Instructions	17

Abstract

...

Beautification of images by Generative Adversarial Networks (GANs)

For millennia, philosophers have been arguing about the nature of beauty. Yet after all these years, we still can't confidently claim to understand what beauty even *is*. Furthermore, the fundamental question whether beauty is subjective or objective remains unsettled (Sartwell, 2022). In this paper, we will use state-of-the-art machine learning techniques to try and create a so-called 'visual definition' of beauty.

Defining Beauty

Throughout history, there have been different conceptions of beauty. The classical conception of beauty, often embodied in classical art, typically concerns the arrangement of independent parts into a coherent whole through properties such as symmetry and order (Sartwell, 2022; Wölfflin, 1932). Another important conception is the idealist conception of beauty described by Plato as a perfect unity. This is in contrast to the classical conception which emphasizes individual parts (Sartwell, 2022). 18th century empiricist philosophers on the other hand looked at beauty in a more instrumental, hedonistic sense. Hume (1740) for example argues that beauty exists to give pleasure and satisfaction to the soul.

While most classical philosophers would have argued that beauty is a property of an object that exists outside the mind, Renaissance and later philosophers argue more for a combination of some objective properties with unique subjective appreciations (Sartwell, 2017). The assumption that there are indeed some objective properties to beauty allows us to examine what exactly it means for something to be beautiful.

The shift from a pure objective conception of beauty towards a more nuanced combination of objective properties with subjective perception logically coincided with the start of psychology as an empirical science in the 19th century.

Among these early psychologists a school of thought now known as Gestalt psychology Gestalt psychology relating to classical conception of beauty.

Machine Learning and Beauty

How we appraise a picture is determined by many factors. Many studies have attempted to explore these factors in a large variety of ways (Deng et al., 2017). In this study, we will investigate the factors underlying the aesthetic value of images using a generative adversarial network (GAN). A GAN is essentially a training pair of neural networks in competition with each other, eventually leading to a strong discriminator network which judges the output produced by the generator (Creswell et al., 2018). Goetschalckx et al. (2019) have developed GANalyze, a framework using GANs that allows us to study cognitive properties by means of so-called visual definitions. We will use GANalyze to generate sequences of images that are supposed to possess either *more*, or *less* aesthetic quality. This way, we will be able to examine the underlying factors responsible for making an image either more or less aesthetic by comparing the supposed low-aesthetic to the supposed high-aesthetic images created by GANalyze. We will perform a finer analysis of the aesthetics-related parameters by looking into the hidden layers of the trained GANalyze network.

If we want to learn something about human aesthetics appraisal from a GAN, we must first validate whether the GAN truly ‘understands’ what it means for an image to possess aesthetic quality. To do this, we must validate whether human participants tend to agree with the GAN on which image appears to be more beautiful. We are also interested whether a person’s experience with art and their demographic characteristics affect their proportion of agreement with the GAN. A relation between art experience and agreement with the GAN could have important implications for the validation of the network and the generalizability of the final results.

Methods

Participants

Questionnaire

To measure the aesthetic experience of each participant, we used a shortened version of the Aesthetic Experience Questionnaire (AEQ; Wanzer et al., 2020). We chose this questionnaire based on a number of factors. First of all, the theoretical background of an aesthetic flow experience (Csikszentmihalyi & Robinson, 1990) seems intuitively relevant to our experiment design where participants will be exposed to a large number of images with varying aesthetic qualities. Second, this questionnaire seems fitting because the AEQ is supposed to be very generalizable to all visual aesthetic experiences and not only famous pieces in a museum. Finally, we picked the AEQ because the authors suggest that it should be possible to reduce the items on the six scales and still end up with a sufficient measurement of aesthetic experience.

As the experience one has with art was not the primary focus of our study, we only used the highest loading items from each of the six scales, resulting in a total of six items at the start of our study in order to save time. Each item was presented with a five-point Likert scale, ranging from 0 (*Strongly disagree*) to 4 (*Strongly agree*). To test whether this shortened version of the AEQ was sufficient to explain the latent aesthetic experience, we used confirmatory factor analysis (CFA).

Behavioral Experiment

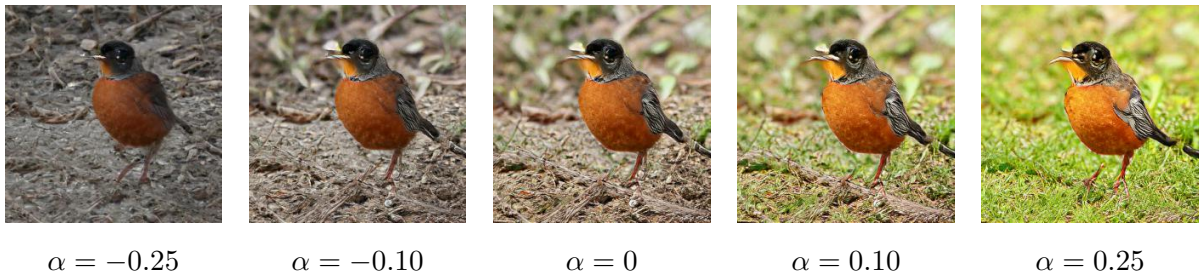
To determine whether GANalyze (Goetschalckx et al., 2019) can be used to study what it means for an image to be aesthetic, I conducted a study to test whether human participants indeed preferred the images that GANalyze generated to be more aesthetic over a corresponding base image.

I used GANalyze with AestheticsNet (Kong et al., 2016) as the assessor and BigGAN-256 (Brock et al., 2019) pretrained on ImageNet (Russakovsky et al., 2015) as the generator to train the GANalyze model for 400,000 iterations. The training resulted in GANalyze being able to produce image sequences based on 1,000 ImageNet categories. For each ImageNet category, I generated three different seeds (i.e., three different images

of Siamese cats) to test for an effect of category. Within each seed, GANalyze produced 21 images with a supposedly increasing amount of aesthetic value (Figure 1). I picked a broad range of α -values with increasing density close to zero to obtain stimuli with obvious changes (e.g., $\alpha = 0.25$), but also subtle changes (e.g., $\alpha = 0.0025$). I chose these values because I hypothesized the participants to agree with GANalyze 100 of the time for trials with a large α -value, and 50 of the time (picking at random) for very small α -values, resulting in a psychometric function. The specific α -values were chosen based on the results of a pilot study with a very broad range.

Figure 1

Truncated Sample of an Image Sequence from One Seed Produced by GANalyze



489 ImageNet categories were manually excluded from the stimulus set because they either contained human faces (which are deliberately warped by BigGAN) or insects such as spiders which might make some participants feel uncomfortable. In addition, stimuli that were unrecognizable were also excluded. This resulted in 9xx categories, each with three seeds and 21 α -values, equaling a total of 30.807 stimuli.

The study was hosted on Prolific.co to obtain a large number of data points in a simple online task, ensuring each of the 30.807 stimuli was evaluated 2-3 times on average. Individuals younger than 17 or those without normal or corrected-to-normal vision were not allowed to participate. At the start of the study, participants were asked to fill in demographic details such as their age, gender, and country. They were also polled on their previous experience and familiarity with visual arts. Afterwards, they had to read and accept an informed consent document after which they were provided with instructions for the task. The participants carried out a spatial two-alternatives forced choice (2AFC) task in which they had to indicate which of the two presented GAN-generated images

Figure 2

Example of a Sequence of Trials in the Study



they considered to be most aesthetically pleasing. They did this by pressing "F" for the left, and "J" for the right image (Figure 2).

Each trial consisted of two images generated from the same seed within an ImageNet category. One of the two images was always the base image with $\alpha = 0$ from that trial's seed while the comparison image was, according to GANalyze at least, either more aesthetic ($\alpha > 0$) or less aesthetic ($\alpha < 0$). The α -value for the comparison image differed for each trial, leading to differences in similarity between the base and comparison images for each trial. An α -value close to zero indicates high similarity and, consequently, difficulty to discriminate between the images. An α -value much larger or smaller than zero indicates low similarity and therefore easy discrimination. The position of the base and comparison image was randomized for each trial to prevent possible confounding factors. Each 10-minute block was created to contain a roughly uniform distribution of chosen α -values and ImageNet categories to ensure that each participant is presented with a comparable stimulus set. Furthermore, the same category never appeared more than one time in a block.

After 10 minutes, the task automatically turned to a screen thanking the participants for their assistance, after which they were given monetary compensation of £1.38 converted to their local currency.

Data Analysis

Structural Equation Modeling

The first requirement for the model is that the latent variable measuring aesthetic experience was properly explained. In addition, we wanted to know whether age, sex, and nationality affect this. Furthermore, we are interested in the effects of aesthetic experience, age, sex, and nationality on the outcome variable of the behavioral task measured by percentage of agreeing with the neural network. To test this, we used structural equation modeling (SEM) with the `lavaan` package (Rosseel, 2012) in R version 4.1.0 (R Core Team, 2021) to build and evaluate models based on theory that would give an explanation to the data. All models were fitted using a maximum likelihood estimation method and the fits were evaluated with the comparative fit index (CFI), Tucker-Lewis index (TLI), and the root mean square error of approximation (RMSEA).

First of all, we performed a priori CFA to validate whether our shortened version of the AEQ was indeed sufficient in explaining the latent aesthetic experience variable. We used the ‘emotional’, ‘cultural’, ‘perceptual’, ‘understanding’, ‘flow-proximal’, and ‘flow-experience’ scales as indicators for our latent variable.

Next, we moved on to more complicated structural models to better describe the full dataset. More specifically, we created a multiple indicators multiple causes (MIMIC) model aimed to show whether age, sex, and nationality affected the latent aesthetic experience variable. Crucially, we were interested in the effects of the aesthetic experience latent variable and age, sex, and nationality on the outcome of the image comparison behavioral task. We used the proportion of agreement with the GAN as the measure for aesthetic judgement here. Our initial idea for this model consisted of the latent variable, explained by the indicators, affecting aesthetic judgement in the behavioral task. Age, sex and nationality would have an effect on both the latent variable and the behavioral outcome of judgement. We also made another alternative model based on different theoretical arguments. For this second model, we decided that it did not make sense for nationality to have a direct effect on aesthetic judgement, so we removed that effect and instead added an effect of nationality on the cultural scale from the AEQ, as

we believed this would make sense. These models would additionally test the questions Wanzer et al. (2020) had regarding their homogeneous population sample.

Even though the new model MIMIC did not appear to be significantly better than the older one, we decided to keep the newer model as it made more theoretical sense.

Behavioral Task

Outliers, total participants, etc.

(Merge with SEM perhaps)

Table 1: Model Fit Indices

Model	χ^2	<i>df</i>	CFI	TLI	RMSEA
CFA	33.92***	9	.96	.93	.07
MIMIC 1	68.72***	29	.94	.91	.05
MIMIC 2	68.00***	29	.94	.91	.05

*** $p < .001$.

Results

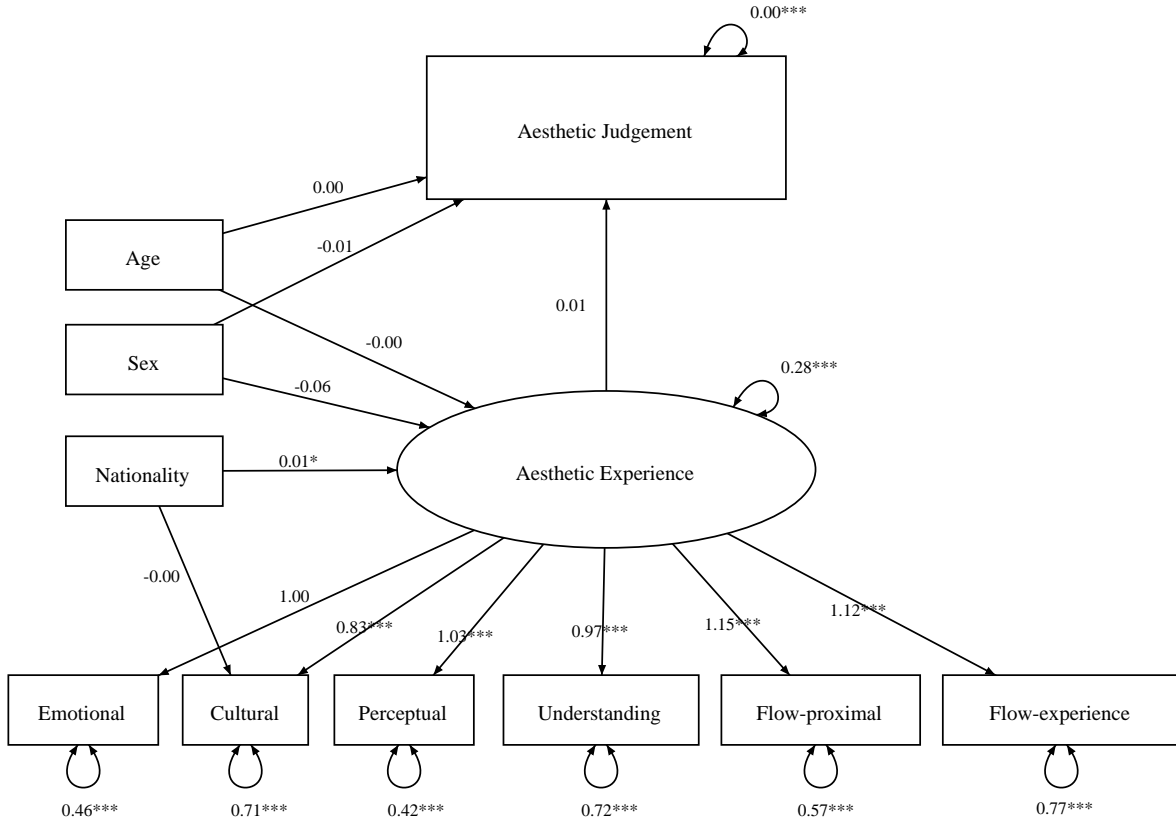
Structural Models

Our first model, the CFA on the shortened AEQ, resulted in a validation of the shortened questionnaire. The model fit can be seen in Table 1 under CFA. All fit indices, apart from the conservative χ^2 , were in the acceptable range (i.e., $> .90$ for CFI and TLI, and $< .08$ for RMSEA). All standardized factor loadings from the latent variable on the measured indices were significant ($p < .001$) and strong enough to justify interpretation ($> .80$). We interpret these results as evidence for our shortened version of the AEQ being adequate for measuring an individual’s aesthetic experience.

For our first MIMIC model, we also obtained good fit measures, as seen in Table 1 under MIMIC 1. This model included an effect of sex, age, and nationality on aesthetic experience, and an effect of sex, age, and nationality on the proportion of agreement with the GAN. While the measurement part of the model remained largely unchanged, the new structural part of the model only showed a significant effect of nationality on aesthetic experience ($\beta = 0.01, p < .05$). This finding may have consequences for the AEQ itself, as Wanzer et al. (2020) had mentioned that their own sample only included US participants. On the other hand, we did not find an effect of age and sex on aesthetic experience, unlike Wanzer et al. (2020). The most important finding here is that there was no effect of the latent aesthetic experience on the behavioral outcome variable of aesthetic judgement, which has important implications for the rest of our study.

For our second MIMIC model, illustrated in Figure 3, we made changes on theoretical bases. The fit indices were equally good, as seen in Table 1 under MIMIC 2, but we believe this second model makes more sense from a theoretical point of view. The standardized factor loadings were again very similar. Again, the effects of age and sex on aesthetic

Figure 3
Final MIMIC Model



experience and aesthetic judgement were not significant. Nationality still has a significant effect on aesthetic experience here, but not on the cultural indicator. There was also no significant effect of aesthetic experience on aesthetic judgement here.

We can interpret these results as there being no effect of a person's demographic characteristics and their aesthetic experience on the outcome of the behavioral image rating task. This implies that the simple image comparison we used in our experiment is not influenced by seemingly relevant factors, pointing to a potential universal nature of aesthetic appreciation of the images we provided in the task. Furthermore, this will make our interpretations of the factors determining the aesthetic value of an image through GANalyze more generalizable.

Behavioral Task

Descriptive

Psychometric Function

To test the effect of different α -values on image preference, we fitted a psychometric function on the data (Figure 4a). The figure shows that the relation between the α -value used for image generation and the behavioral responses can be reasonably fitted with the typical psychometric function commonly found in psychophysical discrimination tasks. As expected, the estimated point of subjective equality (PSE) for the neutral and non-neutral image is at $\alpha = 0$ (99% CI $[-0.00, 0.00]$). The estimated guess rate (equivalent to the lapse rate of negative α 's) is 0.07 (99% CI $[0.07, 0.08]$) and the estimated lapse rate is 0.17 (99% CI $[0.17, 0.18]$).

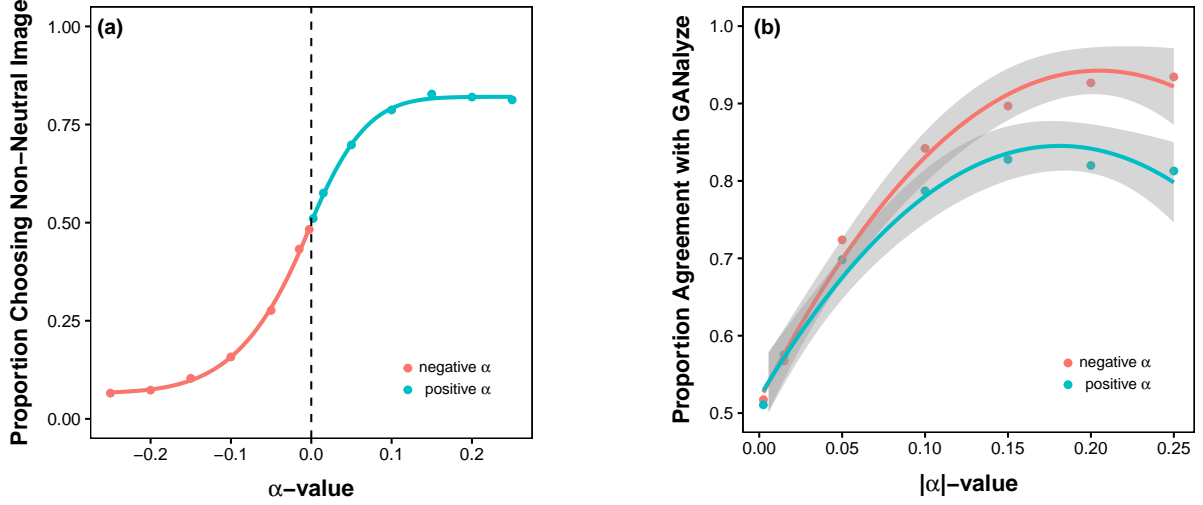
We found a difference between the rating of high- $|\alpha|$ images for positive and negative α -values (Figure 4b). People tend to agree more with GANalyze when it generates images that are supposed to be of lower aesthetic value. This difference is especially pronounced for the extreme variations. A second order polynomial multiple regression with the predictors $|\alpha|$, $|\alpha|^2$, and dummy variable `positive` explained 95% of the observed variance ($R^2 = .95$, $F(3, 10) = 90.59$, $p < .001$). $|\alpha|$ and $|\alpha|^2$ both significantly predicted the proportion agreement with GANalyze ($\beta = 3.87$, $p < .001$; $\beta = -10.03$, $p < .001$). The `positive` variable was also significant ($\beta = -0.05$, $p < .05$), indicating an asymmetry in the appreciation of the generated high- and low-aesthetic images.

Confounding Factors

There was no effect of broad category (), no effect of image position (), no effect of .

Figure 4

Behavioral Results from the Image Rating Task



(a) Psychometric function for the proportion of chosen non-neutral images, illustrating that the estimated threshold value is approximately 0. Choosing the non-neutral image for a negative α -trial corresponds with a disagreement with GANalyze and vice versa. (b) Proportion of agreement with GANalyze with positive and negative α 's taken together, fitted with second order polynomial multiple regression.

Discussion

Behavioral Experiment

Good psychophysical relation between α -value and rating. Positive slope shows that participants are in agreement with GANalyze, therefore we argue that using GANalyze to study the aesthetic properties of images is justified.

Steep curve represents good stimulus representation

GANalyze Parameters

st2

References

- Brock, A., Donahue, J., & Simonyan, K. (2019, February 25). *Large Scale GAN Training for High Fidelity Natural Image Synthesis*. Retrieved May 17, 2021, from <http://arxiv.org/abs/1809.11096>
- Creswell, A., White, T., Dumoulin, V., Arulkumaran, K., Sengupta, B., & Bharath, A. A. (2018). Generative Adversarial Networks: An Overview. *IEEE Signal Processing Magazine*, 35(1), 53–65. <https://doi.org/10.1109/MSP.2017.2765202>
- Csikszentmihalyi, M., & Robinson, R. E. (1990). *The art of seeing: An interpretation of the aesthetic encounter*. Getty Publications.
- Deng, Y., Loy, C. C., & Tang, X. (2017). Image aesthetic assessment: An experimental survey. *IEEE Signal Processing Magazine*, 34(4), 80–106.
- Goetschalckx, L., Andonian, A., Oliva, A., & Isola, P. (2019). GANalyze: Toward Visual Definitions of Cognitive Image Properties. *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 10. <https://doi.org/10.1109/ICCV.2019.00584>
- Hume, D. (1740). *A treatise of human nature*. Courier Corporation.
- Kong, S., Shen, X., Lin, Z., Mech, R., & Fowlkes, C. (2016). Photo Aesthetics Ranking Network with Attributes and Content Adaptation. In B. Leibe, J. Matas, N. Sebe, & M. Welling (Eds.), *Computer Vision – ECCV 2016* (pp. 662–679). Springer International Publishing. https://doi.org/10.1007/978-3-319-46448-0_40
- R Core Team. (2021). *R: A language and environment for statistical computing*. R Foundation for Statistical Computing. Vienna, Austria. <https://www.R-project.org/>
- Rosseel, Y. (2012). Lavaan: An r package for structural equation modeling. *Journal of statistical software*, 48, 1–36.
- Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A. C., & Fei-Fei, L. (2015). ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision*, 115(3), 211–252. <https://doi.org/10.1007/s11263-015-0816-y>

- Sartwell, C. (2017). *Entanglements: A system of philosophy*. SUNY Press. <https://books.google.be/books?id=qDUrvgAACAAJ>
- Sartwell, C. (2022). Beauty. In E. N. Zalta (Ed.), *The Stanford encyclopedia of philosophy* (Summer 2022). Metaphysics Research Lab, Stanford University.
- Wanzer, D. L., Finley, K. P., Zarian, S., & Cortez, N. (2020). Experiencing flow while viewing art: Development of the Aesthetic Experience Questionnaire. *Psychology of Aesthetics, Creativity, and the Arts*, 14(1), 113–124. <https://doi.org/10.1037/aca0000203>
- Wölfflin, H. (1932). Principles of art history, trans. *MD Hottinger (New York, 1932)*, 229.

Appendix

Appendix A: Shortened Aesthetic Experience Questionnaire

In general, when I view art...

1. I experience a wide range of emotions. (emotional)
2. I compare the past culture of the art with present-day culture. (cultural)
3. The composition of a work of art is important to me. (perceptual)
4. I try to understand the work completely. (understanding)
5. I have a clear idea of what to look for when viewing the work of art. (flow-proximal)
6. I lose track of time when I view the work of art. (flow-experience)

Appendix B: Task Instructions