

A Work Project, presented as part of the requirements for the Award of a Master's degree in
Finance from the Nova School of Business and Economics.

MODERN ART MEETS MACHINE LEARNING: ADVANCING PRICE PREDICTIONS
WITH VISUAL FEATURES

NIKLAS ZENG MÀRQUEZ

Work project carried out under the supervision of:

Miguel Lebre Freitas & Luis Catela

15-12-2024

Abstract

This thesis explores machine learning's (ML) potential in art price prediction by integrating traditional artwork features with visual data from high-quality images. Using a dataset of modernist artworks auctioned by Christie's (2007-2024), the study employs hedonic regression, advanced ML, and Convolutional Neural Networks. Results show ML models outperform traditional methods, revealing biases in auction house estimates. Visual features are relevant predictors of price and improve predictions modestly. Auction house valuation biases when combined with ML modes can be interpreted and subsequently mitigated. These findings highlight the promise of data-driven, scalable art valuation frameworks, advancing both academic research and industry practice.

Keywords: Alternative Finance, Art Valuation, Modern Art, Machine Learning, Forecasting, Regression, Convolutional Neural Network

This work used infrastructure and resources funded by Fundação para a Ciência e a Tecnologia (UID/ECO/00124/2013, UID/ECO/00124/2019 and Social Sciences DataLab, Project 22209), POR Lisboa (LISBOA-01-0145-FEDER-007722 and Social Sciences DataLab, Project 22209) and POR Norte (Social Sciences DataLab, Project 22209).

1 Introduction

In the world of art, the story of Vincent van Gogh's "The Red Vineyard" stands out as a poignant example of the complexities in art valuation. Sold during Van Gogh's lifetime for just 400 francs, this piece remains the only work known to have been sold by the artist before his death. Today, Van Gogh's paintings are considered to be some of the most expensive in the world, with pieces such as "Orchard with Cypresses" fetching \$117 million in 2022. This stark disparity highlights the fundamental challenges of art pricing and valuation.

Auction houses, the primary venues for high-value art transactions, employ expert evaluators who assess each piece through a combination of historical research, market knowledge, and comparative analysis. However, these traditional methods face significant limitations: they are labor-intensive, only available before sale, and susceptible to human biases and inconsistencies.

In this thesis we address these challenges by investigating art price prediction using machine learning. Specifically, the research focuses on three critical objectives: First, to assess the performance of machine learning models in predicting art prices; second, to explore the potential importance of visual features in art valuation by extracting and analyzing characteristics through convolutional neural networks (CNNs); and third, to compare machine learning predictions with auction house estimates and investigate the potential of ML to mitigate systematic estimation biases.

The approach builds upon existing literature that has predominantly relied on linear hedonic regression models, which fail to capture the complex, nonlinear relationships that characterize art valuation (Edwards 2004; Aschenfelter and Graddy 2003). Machine learning offers a transformative opportunity, with the potential to leverage high-dimensional data and identify nuanced interactions that traditional methods might overlook. By integrating traditional artwork characteristics and visual features, we add to the literature that visual characteristics

matter in estimating art prices and that machine learning can be leveraged to mitigate auction house biases.

The following sections will provide a comprehensive framework for this study. First, an examination of the pertinent literature surrounding the art market, art price estimation, and machine learning will lay the foundation. This is followed by an overview of the data and the employed methodology. Finally, the findings are presented and discussed, and the study's limitations are acknowledged.

2 Literature Review

Most societies in human history, no matter how developed their material standard of living, have shown creative and artistic endeavors, often highly valued by their respective societies (Farthing 2017). While artistic works are a staple in human history, varying cultural, geographic, and historical factors have influenced radically different creative expression across societies. This heterogeneity is an evident characteristic of art (Adajian 2024), as its creators utilize different techniques and media, appeal to different human senses, and express varying conceptual ideas. This diversity of art has been a significant challenge in establishing consistent valuation methods.

2.1 The Art Market

In contrast to liquid markets such as the stock exchange, where each share of a company is homogenous and trades for the same value, artworks tend to be inherently unique (Aschenfelder and Graddy 2003). Uniqueness and limited availability, where ownership may not change for decades, results in infrequent and time-consuming transactions (Mei and Moses 2002; Lovo and Spaenjers 2018). Given these complexities, auctions have emerged as the dominant venue for high-value art transactions. Auction systems allow individuals to express their subjective valuations through competitive bidding processes. Typically, these auctions operate using an "ascending price" format, where bidding incrementally rises until a final bid is accepted. The

resulting "hammer price" can then be thought to represent the fair market value of a piece (Ashenfelter and Graddy 2003). Nevertheless, auctions have their strategic intricacies. Sellers may set reserve prices, often kept confidential, below which an artwork will not be sold. When bidding fails to meet this reserve, the piece remains unsold. Due to this dynamic, the role of auctioneers, which manage bidding activity and ensure market engagement, becomes especially relevant.

2.2 Auction House Price Estimates

A common practice among auction houses is to set a low and high price estimate before the sale of an artwork. These estimates are expert evaluations for each item, considering the piece's historical significance, condition, prices achieved on similar pieces, and market knowledge. As a result, these estimates provide essential information to both buyers and sellers: buyers are given an idea of the item's potential value, helping them formulate their bids, while sellers receive a range of possible payoffs, which can guide them in setting a reserve price. In the related market of wine auctions Ashenfelter (1989) argued that auction house estimates are generally truthful, showing a high correlation between wine prices and auction estimates. John Abowd and Ashenfelter (1988) further supported this by demonstrating that auction estimates are significantly better price predictors than hedonic regression models.

However, while many studies affirm the predictive accuracy of auction house estimates, systematic under- and overpredictions still occur. For contemporary artworks sold between 1980 and 1994 Beggs and Graddy (1997) find that larger works are generally undervalued while younger paintings are overvalued. Luc Bauwens and Ginsburgh (2000) found similar discrepancies in their study of 1,600 sales of English silver at Christie's and Sotheby's, revealing that Christie's tended to underestimate, whereas Sotheby's overvalued less expensive items and undervalued more costly ones. These systematic prediction errors raise the question whether auction houses pursue a strategic plan through their estimate setting. Beyond potential

strategic biases, Beggs and Graddy (2009) highlighted cognitive biases in the art market, showing that bidders are susceptible to anchoring, whereby auction house estimates and past prices strongly influence "hammer prices."

2.3 Price Determinants of Artworks

While art experts consider various factors when estimating a piece's price, one might assume that intrinsic properties, such as its aesthetics, quality and appeal to viewers, could be the primary determinants. However, examples such as Van Gogh's "The Red Vineyard" show that paintings once valued as 'cheap' can significantly increase in value. Furthermore, artworks identified as forgeries often lose nearly all their value despite no physical alteration to the work (Bailey 2020). This highlights the significant role of factors external to the artwork's visual and material characteristics, such as the artist's reputation at the time of sale.

Mastandrea and Crano (2019) experimented with 309 non-expert participants, showing them identical artworks. Their results revealed that works linked to famous artists were perceived as more interesting and beautiful than identical pieces attributed to non-famous creators. Participants were also willing to pay more to see works they believed to be by well-known artists, illustrating how artist status impacts perceived value. Similarly, in the world of musical art, Lynn, Walker, and Peterson (2016) demonstrated that popularity enhances perceived likeability for lower-quality songs, though this effect does not consistently apply to high-quality works. This evidence suggests that status, reputation, and familiarity are potent determinants of perceived quality and value in art.

In non-experimental settings, the hedonic regression model is frequently applied to explain art prices through various factors such as the artist, genre, technique, dimensions, and the piece's market history (Ginsburgh et al. 2006). This method assumes that the price is a linear function of the artwork's characteristics, where the coefficients are interpreted as the willingness to pay for the underlying feature of the painting (Worthington and Higgs 2006). Typical findings

using hedonic regression indicate that larger artworks and those signed and dated tend to be valued higher, likely due to physical presence and authenticity signals (Aubry et al. 2022; Anderson 1974). The latter was further stressed by Pierre Etienne, Vice President and Deputy Chairman at Christies, who provided us with insights into the price estimation process of auction houses. While the artist, the craftsmanship and similarity to previously sold paintings greatly matter during price estimation, another important factor is the context of time and current consumer demands. As he showcased, a portrait of Toussaint Louverture, a black Haitian general who led the Haitian revolution and slave salvation, can be valued around \$200,000 to \$400,000 in 2024 while in 2014 it would have fetched a mere \$20,000, due to the current importance of equal rights and anti-racism movements¹.

2.4 Machine Learning for Asset Price Prediction

The manual appraisal of artworks is inherently slow and labor-intensive (Bailey 2020). Machine learning presents an opportunity to overcome these limitations by offering more scalable, faster, and potentially more frequent valuation processes (Bailey 2020). While not necessarily more accurate than human experts on an individual case basis, ML-driven models can provide continuous and widespread market transparency (Lang and Maffett 2011; Pagano and Röell 1996).

ML models excel because they can leverage high-dimensional data, perform variable selection through regularization, and identify complex nonlinear interactions among predictors (Gu, Kelly & Xiu 2020). ML approaches demonstrate improved predictive capabilities compared to traditional hedonic regression models (Aubry et al. 2022). It also allows for the integration of less tangible factors, such as social media presence and descriptive text analysis, which Powell et al. (2019) found to correlate with sales outcomes.

¹ Pierre Etienne (Vice President & Deputy Chairman, Christies), in discussion with expert at Christie's, Paris, November 16, 2024

Convolutional Neural Networks trained solely on visual aspects of artworks fail to predict prices better than random chance, as demonstrated by Ayub, Orban and Mukund (2017). This suggests that visual information alone, without contextual data like the artist name, may not capture the intricacies of art valuation, underlining previous findings on the minimal importance of visual attributes (Mastandrea and Crano 2019). Similarly, findings in related fields, such as real estate valuation using computer vision, indicate that while images provide some predictive power, they are typically overshadowed by other intrinsic variables (Glaeser, Kincaid and Naik 2018). The latter findings motivate our exploration of visual characteristics' impact on model performance if combined with traditional artwork features, similar to Aubrey et al. (2022).

Utilizing ML models such as Extreme Gradient Boosting or Random Forests, we will explore how models that are able to capture complex interactions perform in predicting art prices. Furthermore, we will assess the importance of visual characteristics regarding an artwork's valuation by extracting visual features from artwork images through a pre-trained CNN. Based on the previous literature, this should have little significance (Ayub, Orban, and Mukund 2017; Mastandrea and Crano 2019; Glaeser, Kincaid and Naik 2018). Lastly, the auction house estimates are compared against the ML predictions. While we expect auction house estimates to be the best price predictor, the absence of systematic biases is not guaranteed (Beggs and Graddy 1997; Luc Bauwens and Ginsburgh 2000). Hence, utilizing ML to understand those potential biases will be of interest.

3 Data

The data for this study was obtained by scraping sold lots from Christie's using a custom-built Python tool, focusing on 44 famous modernist artists. Typically defined as works created between the 1860s and 1970s, modern art is characterized by experimentation and a break from traditional styles, making it a compelling segment for analysis due to its historical and

commercial impact (Ashenfelter and Graddy 2006). Our sample predominantly includes artists active in this era, except for Jean-Michel Basquiat, who was born in 1960 and created most of his works in his twenties.

After scraping, the dataset initially comprised 8,540 observations from 2007 to the 26th of September 2024. To ensure data robustness, several filtering and cleaning steps were applied, reducing the dataset to 7,755 observations. This involved removing artists with fewer than 35 observations, leaving 37 artists, and excluding non-unique works such as prints and multiples, as well as sculptures, ceramics, and jewelry. Entries lacking realized prices or price estimates were also omitted.

The dataset contains three groups of features that can be summarized as (1) **artwork**, (2) **visual**, and (3) **year** of sale. The artwork set includes direct intrinsic properties of the artwork, such as *dimensions* or *artist*. The visual data represents the visual characteristics of each artwork, extracted from high-quality images of the sold pieces. The exclusion of external features, such as auction location, accommodates the hypothesis that auction house price estimates predominantly reflect the value of intrinsic properties of the artwork at the current time. The latter justifies the inclusion of the year of sale variable, which captures different consumer preferences and demands across time, which play an essential role, as explained by Pierre Etienne². An overview of all features and their descriptions can be found in Appendix 1. Price variables were converted to USD using the spot exchange rate at the time of sale and adjusted to real prices using the January 2007 US CPI as a baseline.

To capture and analyze media and support characteristics, we developed two variable sets:

- 1 ***Sparse Representation***: Aggregates technique and surface information from descriptions, assigning each artwork one dominant medium and support, based on the

² Pierre Etienne (Vice President & Deputy Chairman, Christies), in discussion with expert at Christie's, Paris, November 16, 2024

first appearance in the lot description. Media categories include 'Painting,' 'Oil Painting,' 'Acrylic Painting,' 'Drawing,' and 'Other,' while supports are 'Canvas,' 'Board,' 'Paper,' and 'Other.' Drawings encompass dry media (e.g., pen, chalk) focused on lines, while paintings utilize wet media (e.g., watercolor, paint) with brushstrokes or other techniques (Accardi and Benharrouche 2021). Oil and acrylic paintings were categorized separately due to their economic significance.

- 2 **Granular Representation:** Provides a comprehensive, non-exclusive set of techniques and surfaces, capturing mixed usage. This includes 19 techniques and 8 supports, as detailed in Appendix 2.

The two different data representations will provide insights into whether detailed or simplified categorizations yield more accurate predictions facilitating a comparative evaluation of data representations on a consistent dataset.

Table 1:
Descriptive Statistics for Three most Frequent Categories

	N	% of Total	Nominal Hammer Price (\$) in thousands					
			Mean	Min	P25	P50	P75	Max
All	7,755	100%	1,642.4	0.9	37.5	152.7	759.0	88,805.0
Andy Warhol	1,873	24%	833.6	1.9	10.0	46.6	305.0	57,285.0
Pablo Picasso	916	12%	2,322.2	1.4	100.3	331.5	1,708.5	52,030.8
Henri Matisse	590	8%	755.2	5.7	32.5	80.5	283.0	45,686.5
Drawing	3,475	45%	293.1	0.9	13.8	41.4	138.9	53,000.0
Oil Painting	1,887	24%	4,635.3	10.7	327.2	1,355.2	4,794.4	88,805.0
Acrylic Painting	1,058	14%	2,027.5	7.4	137.3	409.5	1,440.4	48,843.8
Paper	4,815	62%	428.8	0.9	19.9	64.1	242.5	58,363.8
Canvas	2,252	29%	4,226.9	13.8	283.9	1,031.8	4,353.9	88,805.0
Board	394	5%	2,356.2	5.0	138.4	446.5	1,912.4	51,623.1

Artists in the dataset were born between 1834 and 1960, with artworks created from 1853 to 2007. The mean realized price is \$1,642,443, while the median is \$152,741. Andy Warhol leads the dataset with 1,873 sold lots, averaging a hammer price of \$833,597. The most expensive artwork is Robert Rauschenberg's *'Buffalo II,'* which fetched \$88,805,000 in a 2019 New York auction. Furthermore, we see significant differences in price distribution across different media and supports, where oil paintings and pieces on canvas fetch the highest prices, with a median

of \$1,355,205 and \$1,031,771, respectively. Appendix 3 also shows that the lots sold are not distributed evenly across time, where 2013 experienced the most sales with 723 lots sold. Also, artist distribution is not even, as commercial artists such as Andy Warhol and Pablo Picasso dominate the sample, with 35% of instances being attributed to them (see Appendix 4).

4 Methodology

The following section outlines the methodological approaches that were used in this paper to tackle the three main parts of our analysis.

4.1 Train-Test Split

The dataset is first randomly split into a training set comprising 80% of the observations and a test set with the remaining 20%. The training set is used to estimate the weights and biases, or coefficients in the case of the hedonic regression, that minimize the algorithm’s respective loss function. The learned patterns are then used to make predictions on the “*unseen*” test data. Unlike Aubrey et al. (2022), whose training set included instances up until 2014 and predictions were made for 2015 (time-series split), our approach does not segregate the data based on time. Instead, we adjust for inflation by converting nominal prices to real prices using 2007 as a base year, thereby focusing on the core value drivers of art prices. Additionally, as trend and preference information are available to the auction house for the year of sale and no estimates must be made for the far-reaching future, a non-time-series split is deemed appropriate for capturing the real-time valuation of artworks.

4.3 Preprocessing Pipeline

Prior to model training, we preprocess the data to address missing values and standardize the dataset. Here, two missing values for *dimension* are filled using mean imputation. Categorical variables undergo dummy encoding, with the first category dropped to prevent multicollinearity. All features are then scaled using the following transformation to ensure uniformity across different scales:

$$x_i^{scaled} = \frac{x_i - \min(x)}{\max(x) - \min(x)} \quad (1)$$

This normalization preserves the original distribution of values while enabling fair comparison across features during model interpretation, as all distributions are scaled between 0 and 1.

4.4 Predictive Modeling

Using only the training data, a hedonic regression model is trained that relates the (log) price of a specific artwork to a set of its underlying characteristics, serving as a baseline for subsequent machine learning comparisons. This is done for the two training sets, providing a preliminary insight into the features importances. In both cases an equation of the following type was estimated:

$$p_{i,t} = \alpha + X_i\beta + \gamma_t + \varepsilon_{i,t} \quad (2)$$

where $p_{i,t}$ is the log-transformed real hammer price of painting i , α is the respective intercept, X_i is a vector of features, taking on one of the two previously defined forms (e.g. ‘*sparse*’), and γ_t are year fixed effects. The estimated coefficients using the training sets will subsequently be used to generate price predictions \hat{P}_{HR} on the out-of-sample test set.

Following the hedonic baseline, we train more complex models, including Random Forest (RF) and Extreme Gradient Boosting (XGB), known for their robustness in handling non-linear relationships and interactive effects among features.

RF is an ensemble learning technique that constructs various decision trees during training and returns the average prediction of the individual trees. Each tree in the forest is fitted to a randomly bootstrapped sample of the training dataset. The underlying assumption of RF is that each tree will make different mistakes, so combining the results of multiple trees should increase accuracy.

XGB, on the other hand, leverages the concept of boosting, an ensemble technique that builds trees sequentially, such that each subsequent tree corrects the errors of its predecessors. In XGB, initially the average target value is used as a prediction on all our instances to initialize

the algorithm. Subsequently, new trees are created that predict the residuals of the initialization and then of prior trees. Compounded together it makes the final prediction. Additionally, the algorithm includes regularization terms, which allow it to avoid overfitting to the training data. Overall, these models are suited for the predictive tasks at hand, as they provide increased complexity in estimating non-linear relationships, compared to our hedonic baseline model. Also, unlike neural networks, which also handle complex datasets but require large amounts of data and extensive tuning to achieve optimal performance, RF and XGB provide a more manageable approach without compromising on the ability to model non-linear relationships effectively.

We perform hyperparameter tuning through cross-validation, to make sure our models are generalizable and not overfitted. This process involves systematically adjusting model parameters to find the combination that minimize the prediction error. We conduct cross-validation by dividing the training data into multiple smaller sets, where the model is trained on each subset and validated against the remaining data to ensure robustness and prevent overfitting. R-squared (R^2) will be used to assess the model's performance on the test set, as it quantifies the percentage of variation in our dependent variable that is explained by our model. Additionally, the mean-absolute-error (MAE) will be compared across models to compare the average magnitude of errors in our sets of predictions.

4.5 Visual Features

The most effective model is then further examined by integrating visual features extracted using a pre-trained Convolutional Neural Network, specifically the VGG16 model (Simonyan and Zisserman 2015). CNNs are particularly good at processing data with a grid-like topology, such as images, due to their structured network of layers specifically designed for pattern recognition and feature extraction. By removing the network's final fully connected layer, we derive a 4096-dimensional feature vector that encapsulates the visual diversity of the artworks.

This vector is then condensed via principal component analysis (PCA) to distill the most impactful visual elements, ensuring that the feature set remains manageable and focused, thus avoiding noise from excessive dimensionality. The first few principal components capture the majority of the data's variance, efficiently summarizing the key aspects of visual diversity such as major patterns, colors, and textures, hence simplifying the original dataset while preserving its most informative features. By iterating over varying numbers of PCA components during model training, we finely tune the integration of visual data, seeking to enhance the predictive accuracy of our models without introducing noise. The models' performances are compared using the previously used performance metrics. Lastly, interpretation techniques such as permutation importance³ and partial dependency plots⁴ are applied to understand the effect that our features have on the predictive outcome. This methodology aims to understand the influence of visual properties on an artwork's price.

4.6 Auction House Estimate Comparison

The best performing model is then compared against the average real auction house estimates. Here, the following regression model is estimated, to assess how the different predictions line up with the target values of the test set:

$$p_i = \alpha + \beta \hat{p}_i + \varepsilon_i, \quad (3)$$

where p_i is the log real price of artwork i and \hat{p}_i is the predicted value of i . The slope in tandem with the constant will allow to uncover systematic model biases. Auction house biases and efficiency will then be further explored by including the auction house estimates of the training set into the ML model training, similar to Beggs and Graddy (1997). Permutation importance and partial dependency plots will then be leveraged to understand whether factors apart from

³ **Permutation Importance:** Measures the importance of a ML models' feature, by assessing the increase in prediction error that occurs when shuffling (permuting) the feature's value, while controlling all other variables. Larger increases in prediction error indicate greater feature importance

⁴ **Partial Dependency Plots:** Measures the average effect a feature has on the predicted variable by varying its value while controlling all other features. It shows the average relationship between the feature and the target feature

the auction house estimates explain the price. Here, no features apart from the auction house estimate should show any relevance in determining the price, as all relevant information should have been captured by the auction house estimate.

5 Analysis

5.1 Coefficient Analysis

To preliminarily get an understanding of the employed variables, the coefficients for the hedonic regression conducted on the *sparse* training set will be analyzed (see Appendix 5).

Accounting for year and artist-fixed effects, the *dimensions* of the artwork emerge as the most significant price determinant, exhibiting a coefficient of 6.06 ($p < 0.01$). Regarding media, *oil paintings* show the second highest valuation with a coefficient of -0.15 ($p < 0.05$), compared to *acrylic*, which serves as the reference. The coefficient for *miscellaneous paintings*, which include techniques like watercolor, gouache, or ink wash, is slightly lower than *oil*, while *drawings* are substantially cheaper than all forms of paintings, with a coefficient of -1.30 ($p < 0.01$). The higher price for paintings is logical and was also observed in regression results by scholars such as Edwards (2004). Regarding supports, *canvas* is valued the highest, showing a coefficient of 0.14 ($p < 0.1$) relative to the reference category *board*. The analysis also indicates that *signed* (coef. = 0.57, $p < 0.01$) or *dated* (coef. = 0.17, $p < 0.01$) works command higher prices, supporting the assertion that signatures and dates authenticate and thus enhance the artwork's connection to the artist, aligning with results by previous studies (Edwards 2004; Aubrey 2022). This association is further evidenced by the findings of Mastandrea and Crano (2019) and Anderson (1974), where artworks associated with renowned artists garner higher willingness to pay.

Examining artist-specific effects, *Vincent Van Gogh* significantly enhances price value with a coefficient of 2.25 ($p < 0.01$), contrasted against the reference artist *Alberto Giacometti*. Following *Van Gogh*, *Paul Cezanne*, and *Georges Seurat* fetch the highest prices

when controlling for the other features. Surprisingly, *Andy Warhol's* works show a significant negative differential compared to *Giacometti*, which may be attributed to the high volume of lower-valued pieces, which becomes clear through the large gap between his mean (\$833,597) and median prices (\$46,000), overshadowing his premium works (see Appendix 5).

5.2 Model Comparison

Subsequently, we assess the performance of our models, focusing on the influence of different representations of the media and support characteristics of the artwork.

The XGB model consistently outperforms both the hedonic regression and the RF across both data configurations, establishing its robustness for further study, as it can explain 80.8% and 82.4% of the variance in price for the *sparse* and *granular* test set, respectively. The hedonic regression model predicts surprisingly well on the unseen test data, explaining 69.9% (*sparse*) and 71.3% (*granular*) of the variance in price. The RF model exhibits the weakest performance across all configurations, with an R^2 of 0.635 and 0.634 for the *sparse* and *granular* data test sets, respectively (see Table 2). The underperformance of the RF model, particularly in the *granular* representation, may be attributed to its inherent design, which typically favors classification over regression. By averaging multiple decision trees, RF may dilute the precision required for accurate regression predictions. Conversely, the effectiveness of the hedonic regression model suggests that many price determinants in the art market may align with linear patterns, which this model effectively captures. The XGB model's performance underlines its ability to effectively handle non-linearities and complex interactions among features. Its capacity to adapt and learn from complex patterns provides a significant advantage over the random forest and hedonic regression models.

Table 2:
Model Performance Comparison for *Sparse* & *Granular* Feature Representation

	Sparse			Granular		
	\hat{P}_{HR}	\hat{P}_{RF}	\hat{P}_{XGB}	\hat{P}_{HR}	\hat{P}_{RF}	\hat{P}_{XGB}
<i>R</i>²	0.699	0.635	0.808	0.713	0.634	0.824
<i>MAE</i>	0.904	1.025	0.714	0.892	1.038	0.682

The *granular* representation tends to surpass the *sparse* model, suggesting an information gain from incorporating a broader array of techniques and media, obvious in the best-performing XGB model. Here, R^2 is approximately 2 p.p. higher in the *granular* model than the *sparse* one. This indicates that overlooking secondary techniques may reduce predictive precision. Our *sparse* model approach, which assigns predominant weight to the first listed medium or surface and aggregates different techniques into the predefined umbrella categories such as drawings or paper, contrasts with the *granular* approach that benefits from a broader representation without prioritizing a dominant medium. This comprehensive depiction also contributes to the observed performance improvements in Mean Absolute Error (MAE). Here, we see the same patterns as with R^2 , where the MAE improves for HR and XGB while minimally deteriorating for the RF. Lastly, the MAE is noticeably lower for our XGB models compared to our baseline hedonic regression. For instance, our top-performing XGB model's MAE indicates an average absolute estimation error of 0.682 log price units, marking a 24% improvement (0.21 units lower) over the *granular* hedonic model.

These findings underscore the effectiveness of the XGB model in capturing the complex dynamics of art valuation, which is reflected in its superior predictive performance. Lastly, we find that including a non-mutually exclusive representation of the techniques employed in the artwork provides an information gain, which, on average, increases predictive performance.

5.3 Visual Features

This part will evaluate the inclusion of visual features on the best performing XGB model trained on the *granular* feature set.

The dimensionality of the 4096 features extracted from the CNN was reduced to 300 visual principal components using PCA, explaining 84% of the variance in the visual characteristics of our sample. To validate these extracted features' meaningfulness, we calculated the cosine distances between visual feature vectors of randomly selected images and all other artworks. This analysis revealed that artworks deemed similar within the visual dimension share common traits, as illustrated in Appendix 6, where each row depicts five similar works. The visual components effectively capture significant aspects of paintings, discerning color similarities, differentiating between motifs (e.g., landscapes vs. portraits), and identifying distinct styles and stroke patterns.

Through iterative testing of 10, 25, 50, 100, and 300 visual features, it was determined that incorporating 25 visual features yielded the best model performance. This incorporation resulted in modest yet noticeable improvements in performance metrics, with R^2 increasing from 0.824 to 0.830 and MAE decreasing from 0.682 to 0.674 (see Table 3).

Table 3:
Model Performance Comparison for Best Performing XGB with & without Visual Features

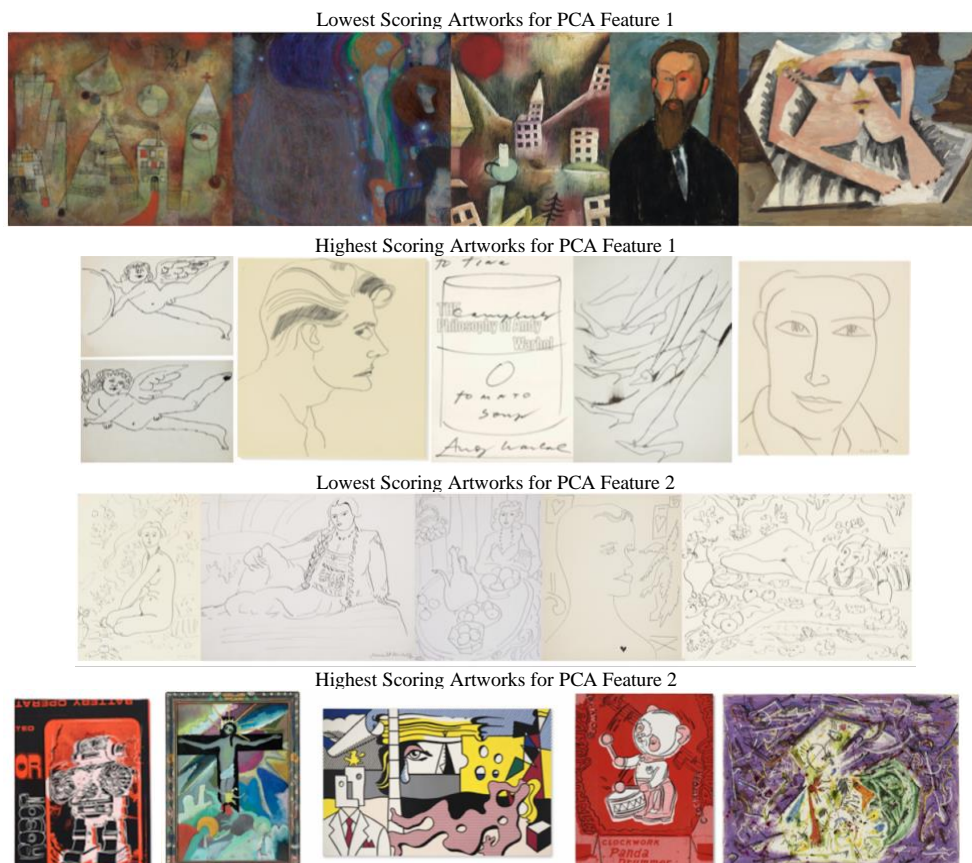
	Granular	
	\hat{P}_{XGB}	$\hat{P}_{XGB+VIS}$
R^2	0.824	0.830
MAE	0.682	0.674

Surprisingly, and contrary to previous findings (Ayub, Orban and Mukund 2017; Mastandrea and Crano 2019), the visual features significantly influence the model predictions, with *PCA Features 1* and *2* emerging as the 4th and 7th most important contributors to price prediction

respectively (see Appendix 7). While *Artist* and *Dimension* remain the predominant determinants of price, the high visual feature importance warrants deeper investigation.

To make sense of the two most influential visual features, we analyzed their highest and lowest scoring artworks (see Figure 1). The 5 works scoring lowest on *PCA Feature 1* are by artists such as Modigliani, Klee, and Picasso, characterized by damp colors, low brightness, little contrast, and expressionist traits with reality presented subjectively, showcasing depth and evoking emotional responses. These works predominantly use *oil on canvas*. Conversely, high *PCA Feature 2* scores correspond to works by artists such as Warhol and Lichtenstein. The feature seems to identify the vibrancy and dynamism within artworks, associating high scores with the utilization of vivid and contrasting colors, innovative forms, and culturally resonant content typical of pop art and more contemporary influences. Scores on the opposite extreme, for both features, are associated with monotone pencil drawings.

Figure 1:
Lowest & Highest Scoring Artworks for *PCA Feature 1* & 2



The partial dependency analysis (see Appendix 8) deepens our understanding of the relationship between features and price. Aligning with the hedonic regression coefficients, *Vincent Van Gogh* leads to the highest prices, while *larger dimensions* or the use of *oil* and *canvas* also on average predict higher realized prices. *PCA Feature 1* exhibits a negative linear trend with price, underlining that damp paintings with expressionist traits are valued higher than drawings. *PCA Feature 2* shows a positive relationship, meaning that, although less impactful than damp paintings, vivid and dynamic colors also lead to higher prices. The economic importance of the two visual features is underlined by quartile analysis (see Table 4): the lowest-scoring 25% of artworks for *PCA Feature 1* achieve a median price of \$588,571, while the highest 25% only reach \$24,137. For *PCA Feature 2*, high-scoring artworks (top 25%) command a median price of \$228,572, compared to \$83,929 for the lowest quartile.

Table 4:
Quartile Price Analysis for *PCA Feature 1 & 2*

	Median Real Price (\$)	
	P25	P75
PCA Feature 1	588,571	24,137
PCA Feature 2	83,929	228,572

The interplay between visual features and traditional art variables reveals a paradox in our analysis. While the incorporation of visual features only marginally improved model performance (R^2 increase 0.006), their high ranking in feature importance, *with PCA Features 1 and 2* emerging as the 4th and 7th most influential predictors, requires a deeper analysis. This discrepancy suggests that the visual features, rather than introducing entirely new predictive information, may entail characteristics already partially encoded in the traditional art features. Our correlation analysis reveals systematic relationships between visual features and specific artistic attributes. *PCA Feature 1*, which we might term "painterly traditionalism", has a moderate negative correlation with oil (-0.51) and canvas (-0.43), indicating that works scoring low in this dimension are more likely to utilize traditional materials, which aligns with our

previous assessment of the low scoring paintings (see Appendix 9). Our previous observations are further underlined by examining the lowest-scoring artists, *Vincent van Gogh*, *Claude Monet*, and *Giorgio de Chirico*, whose works are characterized by atmospheric effects, subtle color gradations, and masterful manipulation of traditional painting techniques (see Appendix 10). Conversely, *PCA Feature 2*, which we could call "compositional dynamism," correlates highest with acrylic (0.24), canvas (0.17), and collage (0.16), reflecting more contemporary approaches. This is evidently captured in the highest-scoring artists, Roy Lichtenstein, Jackson Pollock, and Fernand Léger, whose works embody visual dynamism through bold graphics, action painting, and geometric vibrancy.

The latter demonstrates how the visual features are not independent of other artwork characteristics, showing that certain visual styles can have a higher propensity to encompass certain artists or techniques. As seen above, dynamic and vibrant paintings appear to be painted more often with acrylic paint and by pop art artists. This finding is an explanation for why the performance of the model, integrating the visual features, only marginally increased, while the feature importance of those visual features appeared substantial. Despite the improvement in predictive accuracy only being small, it still is an improvement, justifying the inclusion and importance of integrating the features into modeling. Furthermore, the features can be leveraged to establish a mechanism for quantifying aesthetic characteristics in art valuations.

5.3 Auction House Estimates

Next, regressing the auction house estimates of our test data against the log price reveals that the auction house estimates explain around 95% of the variation in hammer prices. In comparison, the best performing XGB model, which incorporated art features and 25 visual features, scored an R^2 of 83%. Furthermore, the regression of predictions against actuals provides insights into systematic biases. The ML predictions showcase a slope of almost perfectly one and a constant very close to zero, indicating no systematic biases. However, the

auction house estimates have a slope of 0.94 and a constant of 0.96 (see Table 5). This indicates systematic underestimation by the auction house for less valuable artworks. Artworks with a mean estimate around \$22,026 (e^{10}) fetch prices on average 43% higher ($e^{10.36}$), meaning an underestimation of around 30% in that price range.

Table 5:
Regression of Predicted Prices against Realized Prices for Test Data

<i>Model:</i>	XGB: Art + CNN25	AH Estimates	XGB: Art + CNN25 + AH
Prediction	0.997***	0.942***	0.995***
Constant	0.021	0.958***	0.068
<i>N</i>	1551	1551	1551
<i>R</i> ²	0.830	0.949	0.957
<i>MAE</i>	0.674	0.386	0.331

*** $p < .01$, ** $p < .05$, * $p < .10$

The estimation error decreases as average art prices increase. Artworks with a mean real value of \$206,902 ($e^{12.24}$) are underestimated by 21% (e^{12}), and artworks valued around \$15 million and more do not suffer from large estimation errors. This is consistent with findings by Luc Bauwens and Ginsburgh (2000), who also observed systematic underestimation by Christie's auction house.

To assess whether machine learning can mitigate and help explain the systematic bias of the auction house, the XGB model trained on the *granular* data and 25 visual features was enhanced by further conditioning it on the average auction house estimates. The results show that this combined model outperforms the pure auction house estimates. The predictions explain 96% of the variance in the dependent variable, a 1p.p. increase over the sole auction house estimates. Importantly, the systematic bias is entirely mitigated, as the prediction slope equals 1 and the constant is closest to 0. This indicates the ML model is able to account for factors the auction house methodology overlooks.

The observed patterns in feature importance and partial dependence plots underscore this. The importance of features like *Artist*, *Year*, *Dimension*, and *Signed* (see Appendix 12) remains

even when accounting for auction house estimates. If the auction house errors were not systematic, we would not expect these features to maintain predictive power independently.

For example, the partial dependence plot (see Appendix 13) shows that artworks by *Van Gogh* are still predicted to fetch higher prices even when controlling for the auction estimates. This suggests the auction house's estimates do not fully capture the premium associated with certain artists. Economic significance is underlined, as even after holding constant all features including the auction house estimate, *Van Gogh's* paintings are still on average predicted to be 22% ($e^{0.2}$) higher than *Andy Warhol's*.

Additionally, the continued importance of the *Year* feature implies that the auction house estimates are not fully dynamic in adjusting to market trends and shifts in consumer demand over time. Artworks sold in the year 2007 for example are still on average predicted to be 10% ($e^{0.095}$) higher than artworks sold in the year 2020. This points to a systemic bias in the auction house's valuation methodology, which might lead to consistent under- or overestimations in certain periods. The *dimension* feature appears to have an impact, especially on 'very' large artworks, which are predicted to be more expensive, also after accounting for the house estimate. Lastly, the signature, also continues to positively contribute to the price, where signed works are still predicted to be more valuable, with a average price increase of 7% ($e^{0.067}$) over artworks with no signature. Other features like *Age* or certain *PCA_Features* also appear to not be fully accounted for.

The evidence of systematic bias in the auction house estimates underscores the value of integrating machine learning techniques to refine and advance the art valuation process. As demonstrated ML allows to uncover these biases and allow to subsequently mitigate them. Incorporating machine learning into the art pricing landscape could lead to more accurate, dynamic, and comprehensive assessments, benefiting auction houses and the broader art market.

Conclusion

This study sought to explore art price prediction through machine learning. By leveraging predictive models like Extreme Gradient Boosting and integrating both intrinsic artwork characteristics and visual features, the research uncovered insights into art valuation dynamics. A pivotal discovery emerged from the examination of data representation methodologies. While previous studies have presented varying approaches to capturing art media and supports, we demonstrated the significance of a non-mutually exclusive data representation strategy. By allowing for a broader pool of art media and techniques, improved model performance was achieved. However, the methodology remains limited. The two approaches compared represented opposing extremes in data representation—a *sparse* approach using broad umbrella terms and a *granular* approach assigning equal weight to all techniques. It is posited that a hybrid approach, incorporating technique granularity while accounting for technique dominance, could offer the most promising avenue for future investigation.

The most striking discovery emerged from the visual feature analysis. Contrary to previous studies by Ayub, Orban and Mukund (2017) and Mastandrea and Crano (2019), which suggested visual characteristics had no impact on pricing, the research revealed that extracted visual features significantly influenced price prediction. Specifically, PCA-transformed visual features ranked among the top predictors. These features captured sophisticated artistic nuances, such as "painterly traditionalism" and "compositional dynamism," demonstrating machine learning's potential to quantify aesthetic characteristics, which upon manual inspection provide valuable insights into how different styles are valued by the market.

The best-performing XGB model, incorporating *granular* artwork features and 25 visual features greatly outperformed traditional hedonic regression. While the model performed worse than auction house estimates, it notably avoided the systematic biases prevalent in those estimates. The analysis revealed that Christie's systematically underestimates prices of works

worth less than \$15 million, a pattern consistent with earlier observations by Luc Bauwens and Ginsburgh (2000). Reasons for the latter being the systematic inability by the auction house to account for the value of certain artists, time trends, or factors such as the work's dimensions. The methodology's limitations must be acknowledged. As Pierre Etienne highlighted, art valuation is profoundly influenced by contemporary social contexts and consumer demands, dimensions challenging to encode in a machine learning model. As with the portrait of Toussaint Louverture the valuation of artworks can dramatically shift with social movements, a nuance the current approach cannot capture. Furthermore, by not utilizing a time-series split to separate our training and test data, our ML predictions are based not only on information from the past but also the future. This future sales information was logically not available to the auction house when estimating prices, which may limit the directness of our comparative analysis. Lastly, this study remains inconclusive regarding the underlying reasons for the estimation bias observed. It remains unresolved whether Christie's systematically fails to properly incorporate the value of certain artwork characteristics into their valuation, or whether strategic incentives might motivate such observed patterns. The anchoring effect of auction house estimates (Beggs and Graddy 2009), also introduces methodological complexity by potentially influencing the realized price. Consequently, the published estimates may have systematically shaped the final artwork prices, suggesting that the predictive model's performance could be more robust if deployed in a real setting.

Despite these constraints, the predictive performance of the XGB model and its lack of systematic biases suggest significant economic potential. Machine learning models could benefit the art market by providing more scalable, faster, and potentially more frequent valuation processes (Bailey 2020), providing a critical step toward more sophisticated approaches to understanding and predicting art market dynamics.

Bibliography

- Abowd, John, and Orley Aschenfelter. 1988. "Art Auctions: Prices Indices and Sale Rates for Impressionist and Contemporary Pictures." *Mimeo, Department of Economics, Princeton University*.
- Accardi, Angelo, and Yoel Benharrouche. 2021. *Drawing vs Painting*. 18 April. Accessed November 13, 2024. <https://www.eden-gallery.com/news/drawing-vs-painting>.
- Adajian, Thomas. 2024. "The definition of Art." *The Stanford Encyclopedia of Philosophy*.
- Anderson, Robert C. 1974. "Paintings as an Investment." *Economic Inquiry* 12 (1): 13-26.
- Aschenfelter, Orley. 1989. "How Auctions Work for Wine and Art." *Journal of Economic Perspectives* 3 (3): 23-36.
- Aschenfelter, Orley, and Kathryn Graddy. 2003. "Auctions and the Price of Art." *Journal of Economic Literature* 41 (3): 763–786.
- Ashenfelter, Orley, and Kathryn Graddy. 2006. *Handbook of the Economics of Art and Culture*. Edited by Victor A. Ginsburgh and David Throsby. Vol. 1. Amsterdam: Elsevier B.V.
- Aubry, Mathieu, Roman Kraeussl, Gustavo Manso, and Christophe Spaenjers. 2022. "Biased Auctioneers." *Journal of Finance, Forthcoming* 1 - 46.
- Ayub, Rafi, Cedric Orban, and Vidush Mukund. 2017. "Art Appraisal Using Convolutional Neural Networks." *Unpublished* (Stanford University).
- Bailey, Jason. 2020. "Can Machine Learning Predict the Price of Art at Auction?" *Harvard Data Science Review* 2 (2).
- Bauwens, Luc, and Victor Ginsburgh. 2000. "Art Experts and Auctions: Are Pre-sale Estimates Unbiased and Fully Informative?" *Louvain Economic Review*, 66 (2): 131-144.
- Beggs, Alan, and Kathryn Graddy. 1997. "Declining Values and the Afternoon Effect: Evidence from Art Auctions." *The RAND Journal of Economics* 12 (1): 13-26.
- Beggs, Alan, and Kathryn Graddy. 2009. "Anchoring Effects: Evidence from Art Auctions." *The American Economic Review* 99 (3): 1027-1039.
- Edwards, Sebastian. 2004. "The Economics of Latin American Art: Creativity Patterns and Rates of Return." *National Bureau of Economic Research* 10302.
- Farthing, Stephen. 2017. *Kunst. Die ganze Geschichte*. Köln: DuMont Buchverlag.

- Glaeser, Edward L., Michael Scott Kincaid, and Nikhil Naik. 2014. "Computer Vision and Real Estate: Do Looks Matter and Do Incentives Determine Looks." *National Bureau of Economic Research* 25174: 1-36.
- Gu, Shihao, Bryan Kelly, and Cacheng Xiu. 2020. "Empirical Asset Pricing via Machine Learning." *The Review of Financial Studies* 33 (5): 2223-2273.
- Lang, Mark, and Mark Maffett. 2011. "Transparency and liquidity uncertainty in crisis periods." *Journal of Accounting and Economics* 52 (2): 101-125.
- Liu, Amy. 2015. "Art Arbitrage - Violations of the Law of One Price Created by Fine Art Auctions." *Undergraduate Economic Review* 12 (1).
- Lovo, Stefano, and Christophe Spaenjers. 2018. "A Model of Trading in the Art Market." *The American Economic Review* 108 (3): 744-774.
- Lynn, Freda B., Mark H. Walker, and Colin Petersen. 2016. "LynIs popular more likeable? Choice status by intrinsic appeal in an experimental music market." *Social Psychology Quarterly* 79 (2): 168-180.
- Mastandrea, Stefano, and William D. Crano. 2019. "Peripheral Factors Affecting the Evaluation of Artworks." *Empirical Studies of the Arts* 37 (1): 82-91.
- Mei, Jianping, and Michael Moses. 2002. "Art as an Investment and the Underperformance of Masterpieces." *The American Economic Review* 92 (5): 1656-1668.
- Pagano, Marco, and Ailsa Roell. 1996. "Transparency and Liquidity: A Comparison of Auction and Dealer Markets with Informed Trading." *The Journal of Finance* 51 (2): 579-611.
- Pesando, James E. 1993. "Art as an Investment. The Market for Modern Prints." *American Economic Review* 85 (5): 1075-1089.
- Powell, Laurel, Anna Gelich, and Zbigniew W. Ras. 2019. "Developing Artwork Pricing Models for Online Art Sales Using Text Analytics." *International Joint Conference on Rough Sets* 480-494.
- Simonyan, Karen, and Andrew Zisserman. 2015. "Very deep convolutional networks for large-scale image recognition." *Int. Conf. on Learning Representations*.
- Worthington, Andrew C., and Helen Higgs. 2006. "A note on financial risk, return and asset pricing in Australian modern and contemporary art." *Journal of Cultural Economics* 30: 73-84.

Appendix

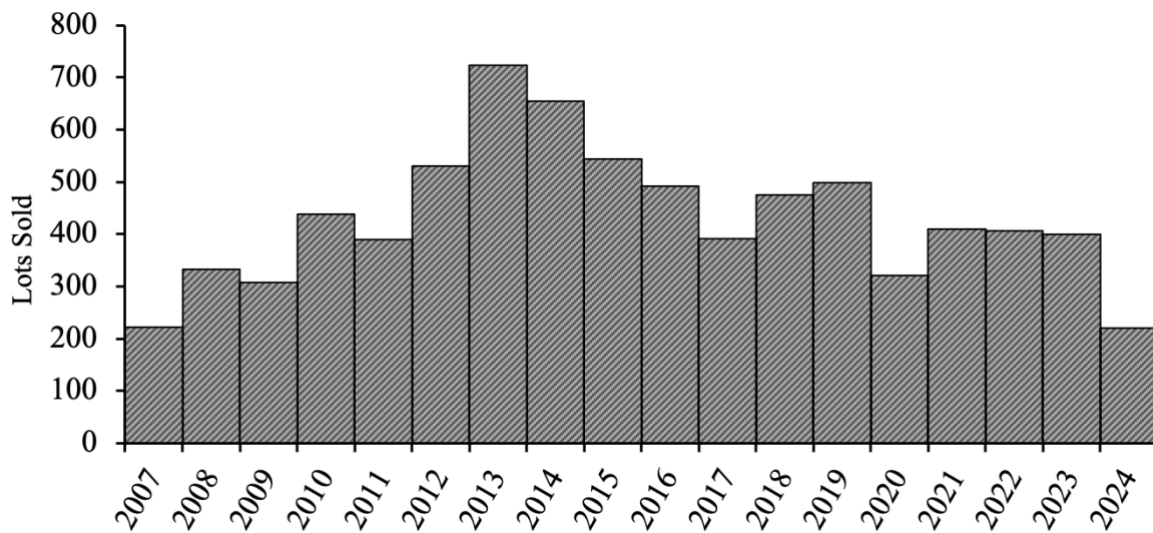
Appendix 1: Variable Overview

Category	Variable	Description
Target	LogPrice	Real hammer price of artwork. Transformed using the natural logarithm. Transformed to real price using monthly US CPI levels, with January 2007 as base
Art	Artist	Categorical variable for artist who created the artwork
	Signed	Binary variable indicating if artwork is signed or not (1/0)
	Dated	Binary variable indicating if artwork is dated or not (1/0)
	LogDimension	Continuous variable representing the surface of the artwork in cm ² . Transformed using the natural logarithm
	Is Alive	Binary variable indicating if artist was alive at the date of sale or not (1/0)
	Age	Continuous variable representing the age of the artist in years at the year of creation
	Medium	Sparse or Granular representation of artwork media
	Support	Sparse or Granular representation of artwork support
Visual	Visual Features	Reduced-dimensional feature vector capturing visual characteristics of artworks (derived from CNN extracted features)
Year	Year	Year of Sale

Appendix 2: Medium & Support Feature Overview

	Medium	Support
Sparse: Mutual Exclusivity	Painting , Oil Painting, Acrylic Painting, Drawing, Other	Canvas , Board, Paper, Other
Granular: No Mutual Exclusivity	Acrylic, Chalk, Charcoal, Collage, Crayon, Gouache, Graphite, Ink, Ink Wash, Mixed Media, Oil, Pastel, Pen, Pencil, Pigment, Sanguine, Tempera, Watercolor, Other	Board, Canvas, Cardboard, Glass, Panel, Paper, Textile, Other

Appendix 3: Sold Lots per Year



Appendix 4:
Artworks per Artist in total Sample (Training + Test Data)

Artist	<i>N</i>	% of Total
Andy Warhol	1873	24.15%
Pablo Picasso	916	11.81%
Henri Matisse	590	7.61%
Salvador Dalí	300	3.87%
Joan Miró	300	3.87%
Willem De Kooning	287	3.70%
Edgar Degas	245	3.16%
Jean-Michel Basquiat	229	2.95%
Giorgio De Chirico	221	2.85%
Paul Klee	221	2.85%
Fernand Léger	195	2.51%
René Magritte	195	2.51%
Marc Chagall	181	2.33%
Claude Monet	165	2.13%
Tamara De Lempicka	160	2.06%
Helen Frankenthaler	135	1.74%
Robert Rauschenberg	135	1.74%
Wassily Kandinsky	128	1.65%
Egon Schiele	126	1.62%
Alberto Giacometti	125	1.61%
Gustav Klimt	118	1.52%
Roy Lichtenstein	116	1.50%
Man Ray	109	1.41%
Paul Gauguin	90	1.16%
Sonia Delaunay	78	1.01%
Georgia O'Keeffe	71	0.92%
Paul Cézanne	60	0.77%
David Smith	50	0.64%
Amedeo Modigliani	48	0.62%
Vincent Van Gogh	48	0.62%
Mark Rothko	47	0.61%
Piet Mondrian	45	0.58%
Georges Seurat	42	0.54%
Francis Bacon	36	0.46%
Jackson Pollock	35	0.45%
Yves Klein	35	0.45%

Appendix 5:
Coefficients of Hedonic Regression on *Sparse Data*

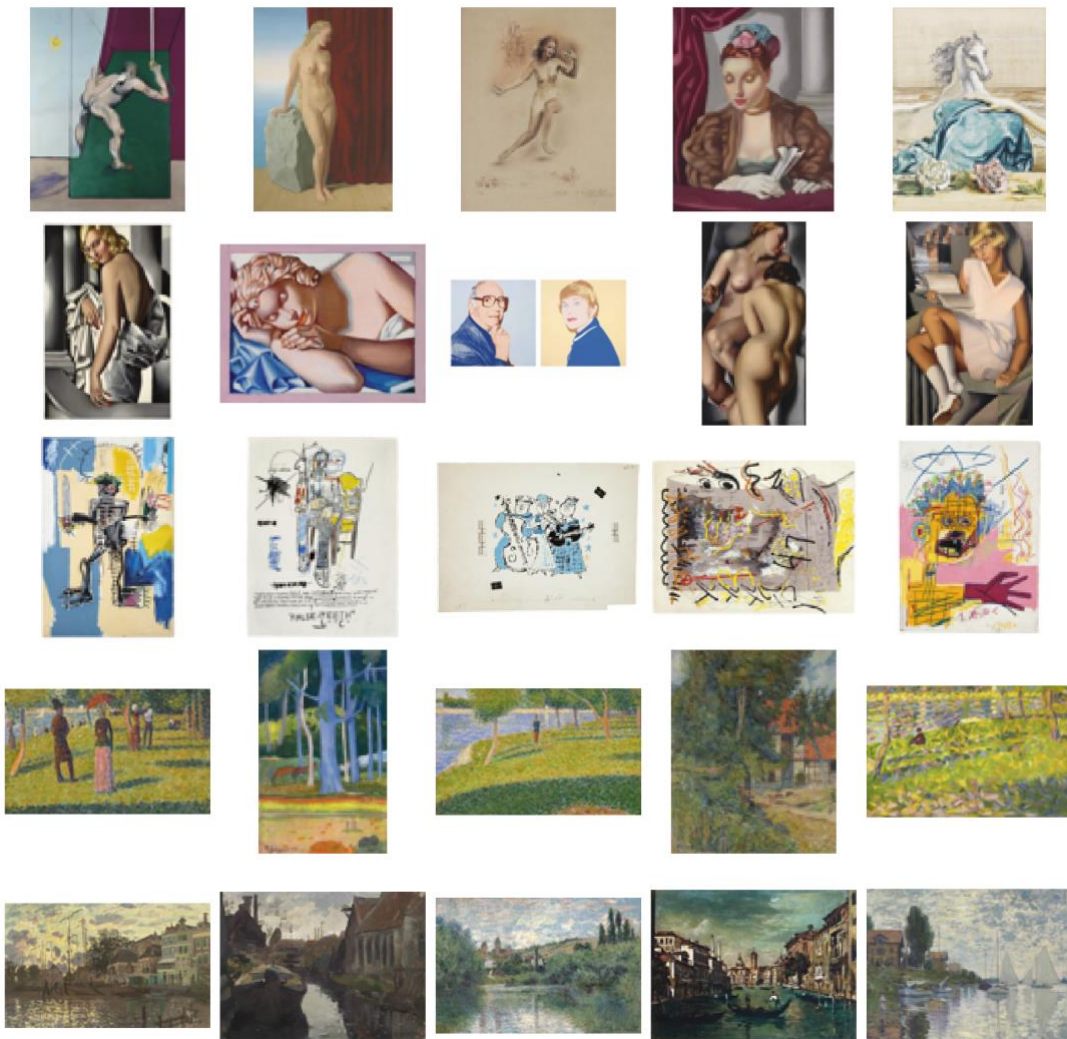
Feature	Coef.
Year FE	YES
Const	11.29***
Log Dimension	6.06***
Signed	0.57***
Dated	0.17***
Age	-0.61***
Is Alive	-0.81***
Medium: Drawing	-1.30***
Medium: Misc Painting	-0.26***
Medium: Oil Painting	-0.15**
Medium: Other	-0.94***
Support: Canvas	0.14*
Support: Other	-1.18***
Support: Paper	-1.22***
Artist: Amedeo Modigliani	0.73***
Artist: Andy Warhol	-1.71***
Artist: Claude Monet	0.52***
Artist: David Smith	-2.09***
Artist: Edgar Degas	0.55***
Artist: Egon Schiele	0.45***
Artist: Fernand Leger	-0.84***
Artist: Francis Bacon	1.17***
Artist: Georges Seurat	1.35***
Artist: Georgia O'Keeffe	0.58***
Artist: Giorgio De Chirico	-1.92***
Artist: Gustav Klimt	-0.54***
Artist: Helen Frankenthaler	-2.63***
Artist: Henri Matisse	-0.16
Artist: Jackson Pollock	0.85***
Artist: Jean-Michel Basquiat	0.28*
Artist: Joan Miro	-0.69***
Artist: Man Ray	-2.04***
Artist: Marc Chagall	-0.19
Artist: Mark Rothko	0.85***
Artist: Pablo Picasso	0.48***
Artist: Paul Cezanne	1.19***
Artist: Paul Gauguin	-0.06
Artist: Paul Klee	-0.25*

Artist: Piet Mondrian	-0.77***
Artist: Rene Magritte	0.47***
Artist: Robert Rauschenberg	-1.38***
Artist: Roy Lichtenstein	-0.30*
Artist: Salvador Dali	-1.09***
Artist: Sonia Delaunay	-2.07***
Artist: Tamara De Lempicka	-1.74***
Artist: Vincent Van Gogh	2.25***
Artist: Wassily Kandinsky	0.33**
Artist: Willem De Kooning	-0.84***
Artist: Yves Klein	0.09

*** $p < .01$, ** $p < .05$, * $p < .10$

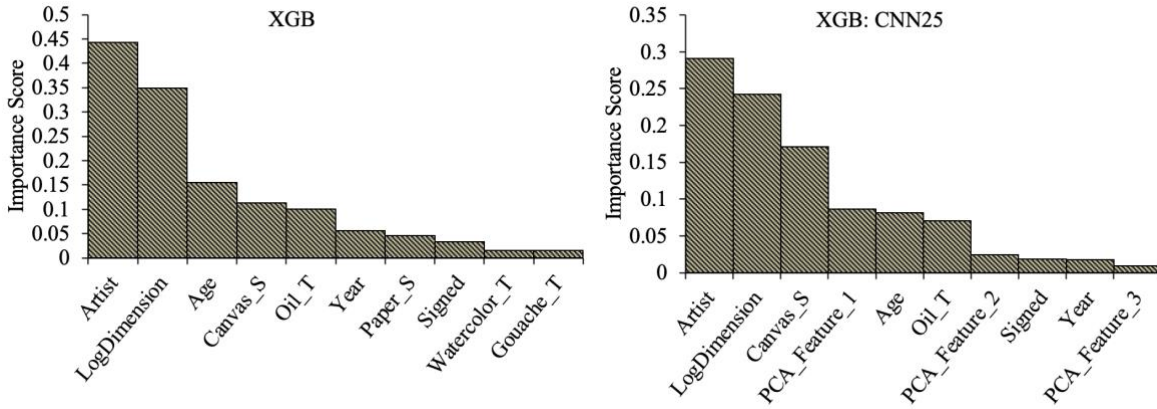
Appendix 6:

Visual Representation of 5 most similar artworks based on Visual Feature Similarity



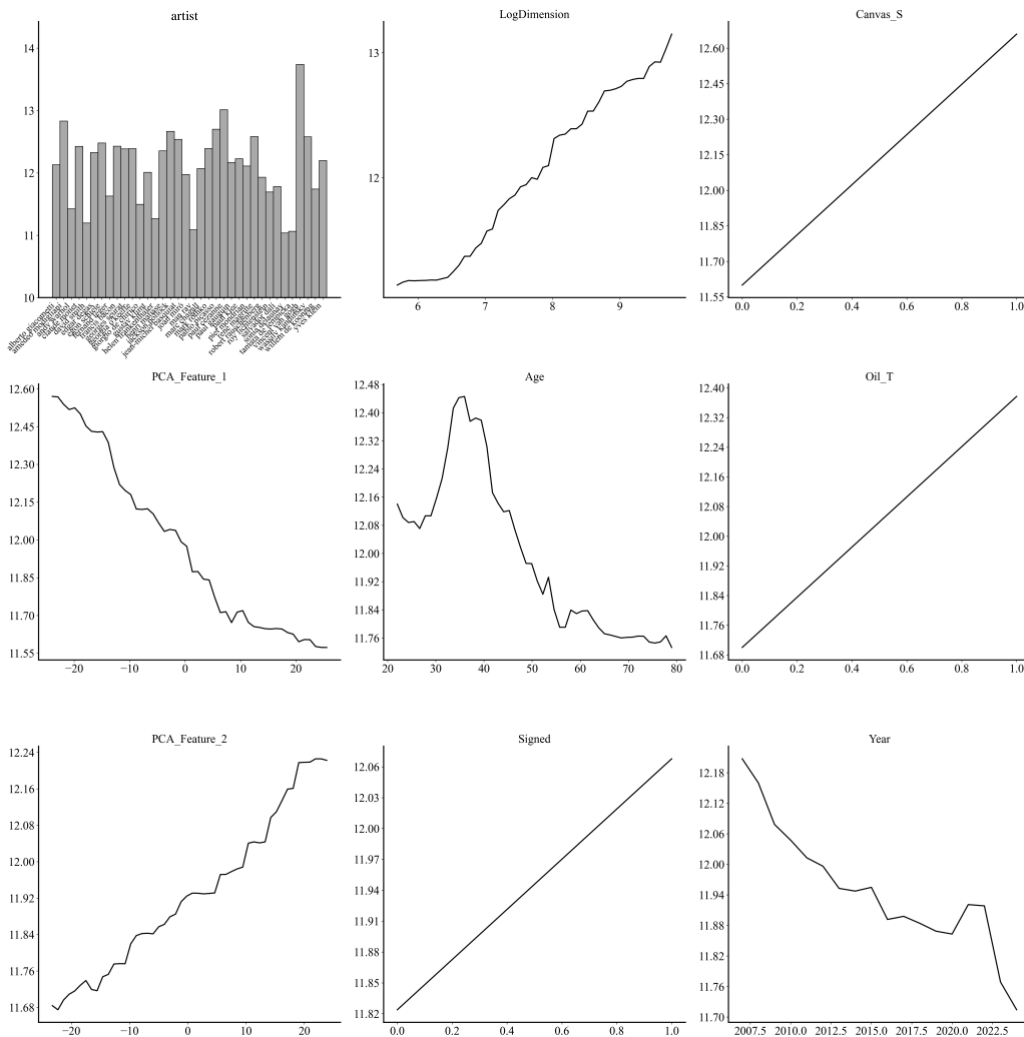
Appendix 7:

Permutation Feature Importance of *Granular* XGB model with and without Visual Features



Appendix 8:

PDP Plot for *Granular* XGB including Visual Features



Appendix 9:
Correlation of PCA Features with Media & Support

	PCA Feature 1		PCA Feature 2	
	Medium/Support	Correlation	Medium/Support	Correlation
Highest	Ink	0.27	Acrylic	0.24
	Pen	0.23	Canvas	0.17
	Graphite	0.2	Collage	0.16
Lowest	Gouache	-0.16	Charcoal	-0.09
	Canvas	-0.43	Pen	-0.15
	Oil	-0.51	Pencil	-0.18

Appendix 10:
Mean PCA Feature Value for Top Three highest & lowest Artists

	PCA Feature 1		PCA Feature 2	
	Artist	Mean	Artist	Mean
Highest	Henri Matisse	8.8	Roy Lichtenstein	14.4
	Andy Warhol	8.5	Jackson Pollock	13.7
	Jean-Michel Basquiat	5.9	Fernand Léger	10.3
Lowest	Vincent Van Gogh	-15.6	Henri Matisse	-10.4
	Claude Monet	-15.8	Claude Monet	-11.5
	Giorgio de Chirico	-16.2	Gustav Klimt	-13.2

Appendix 11:
Regression of Predicted against Realized Prices on Test Data

