Title:

Evaluating XGBoost for Regression on ImageNet-Derived Features Using Statistical Feature Selection

Abstract

Machine learning regression tasks on high-dimensional image data typically demand computationally intensive deep learning models. This study explores an alternative approach by transforming ImageNet images into numerical feature vectors through pre-trained convolutional neural networks, applying statistical feature selection to reduce dimensionality, and utilizing XGBoost for regression modeling. The pipeline's performance is evaluated using Mean Squared Error (MSE), measuring prediction accuracy on synthetic regression targets derived from image metadata attributes. Experimental results demonstrate that statistical feature selection significantly improves computational efficiency and model accuracy, confirming XGBoost's capability to handle high-dimensional, image-derived structured data for regression tasks.

1. Introduction

Image classification and object recognition have been dominated by convolutional neural networks (CNNs) trained on large-scale datasets such as ImageNet. However, with increasing interest in explainable, scalable, and computationally efficient models for tabular and structured data, tree-based ensemble algorithms like XGBoost have garnered significant attention.

This study proposes a hybrid pipeline where image features extracted from ImageNet images via pre-trained deep CNNs are used for regression tasks. Given the high dimensionality of these feature vectors, statistical feature selection techniques are applied to retain only the most relevant features. XGBoost is then trained to predict synthetic continuous targets associated with image metadata, such as average color intensity or image contrast scores.

2. Literature Review

Recent literature highlights the versatility of transfer learning and feature extraction in computer vision (Shin et al., 2016). CNN models such as ResNet and VGG, pre-trained on ImageNet, are often repurposed for downstream tasks including regression and classification on medical images, satellite imagery, and artwork datasets.

Meanwhile, XGBoost (Chen & Guestrin, 2016) has emerged as one of the most powerful and efficient gradient boosting implementations for structured data, regularly outperforming neural networks in tabular machine learning competitions.

Statistical feature selection, including techniques like correlation filtering, ANOVA F-tests, and mutual information analysis, plays a vital role in reducing redundant features, improving model performance, and enhancing interpretability (Guyon & Elisseeff, 2003).

3. Methodology

3.1 Dataset

The ImageNet 2012 dataset includes over 1.2 million images across 1,000 categories. For this study:

- 50,000 images were randomly sampled.
- Synthetic continuous regression targets were generated by computing statistical attributes from image metadata (e.g., mean pixel intensity, standard deviation of color histograms, contrast scores).

3.2 Data Analytics: Statistical Feature Selection

Image features were extracted using a pre-trained ResNet-50 model's penultimate fully connected layer, producing a 2048-dimensional feature vector per image.

Statistical feature selection was applied to reduce dimensionality:

- Variance Thresholding: Removed low-variance features.
- Correlation Analysis: Eliminated highly correlated features (r > 0.9).

 Mutual Information Analysis: Ranked features based on dependency with the continuous target.

The top 300 features (from 2048) were retained for modeling.

3.3 Regression Modeling: XGBoost

XGBoost's gradient boosting decision trees (GBDT) algorithm was employed for regression.

Model Configuration:

• Objective: reg:squarederror

• Number of Estimators: 500

• Learning Rate: 0.05

• Max Depth: 8

• Subsample: 0.8

• Colsample_bytree: 0.7

Hyperparameters were optimized via grid search with 5-fold cross-validation on the training set.

4. Experimental Setup

The experiments were executed on:

• GPU: NVIDIA RTX A5000

• Frameworks: XGBoost 1.7.4, Scikit-learn, PyTorch 2.1

• Data split: 70% Training, 15% Validation, 15% Test

All feature selection and modeling steps were encapsulated in reproducible Python pipelines.

5. Evaluation Metrics

Mean Squared Error (MSE) was the primary metric for evaluating model performance:

$$\label{eq:mse-1} \begin{split} \mathsf{MSE}=&1\mathsf{N}_i=&1\mathsf{N}(\mathsf{yi}-\mathsf{y}^i)2\\ \mathsf{MSE}=&1\mathsf{N}_i=&1\mathsf{N}(\mathsf{yi}-\mathsf{y}^i)2\\ \mathsf{MSE}=&1\mathsf{N}(\mathsf{yi}-\mathsf{y}^i)2\\ \end{split}$$

Lower MSE values indicate better model predictions.

Secondary metrics:

- Root Mean Squared Error (RMSE)
- R-squared (R²)

6. Results and Discussion

6.1 Performance Evaluation

Model Configuration	MSE (Lower Better)	R² Score
XGBoost (No FS, 2048 feats)	0.063	0.712
XGBoost (Stat FS, 300 feats)	0.051	0.775

Key Insights:

• Statistical feature selection improved MSE by 19%, reducing computational cost and training time by approximately 40%.

• Despite dimensionality reduction, model accuracy improved, indicating that most of the original 2048 features were redundant or non-informative.

6.2 Error Analysis

Residual analysis confirmed homoscedasticity (constant error variance) and minimal overfitting, particularly in the feature-selected model. High-importance features were consistently aligned with color-based statistics and texture descriptors.

7. Conclusion

This research demonstrates the feasibility of combining deep CNN-based feature extraction with XGBoost for regression on image-derived structured data. Applying statistical feature selection not only improved computational efficiency but also enhanced model accuracy as measured by MSE.

The findings underscore the adaptability of XGBoost for high-dimensional regression tasks and the value of feature selection in deep feature pipelines, offering a viable alternative to end-to-end deep learning models for certain applications.

8. Future Work

Future research directions:

- Expanding to multi-output regression tasks.
- Integrating embedded feature selection via XGBoost's feature importance scores.
- Comparing against neural network-based regression architectures (MLP, TabNet).
- Applying to real-world use cases: medical imaging scores, satellite imagery regression.

References

- 1. Russakovsky, O., et al. (2015). ImageNet Large Scale Visual Recognition Challenge. *IJCV*.
- 2. Chen, T., & Guestrin, C. (2016). XGBoost: A scalable tree boosting system. KDD.
- 3. Shin, H., et al. (2016). Deep CNNs for Computer-Aided Detection: CNN Architectures, Dataset Characteristics and Transfer Learning. *IEEE TMI*.
- 4. Guyon, I., & Elisseeff, A. (2003). An Introduction to Variable and Feature Selection. *JMLR*.
- 5. He, K., et al. (2016). Deep Residual Learning for Image Recognition. CVPR.