



Munich Personal RePEc Archive

A machine learning approach for assessing labor supply to the online labor market

Fung, Esabella

9 October 2023

Online at <https://mpra.ub.uni-muenchen.de/118844/>
MPRA Paper No. 118844, posted 14 Oct 2023 07:21 UTC

A machine learning approach for assessing labor supply to the online labor market

Esabella Fung

Abstract

The online labor market, comprised of companies such as Upwork, Amazon Mechanical Turk, and their freelancer workforce, has expanded worldwide over the past 15 years and has changed the labor market landscape. Although qualitative studies have been done to identify factors related to the global supply to the online labor market, few data modeling studies have been conducted to quantify the importance of these factors in this area. This study applied tree-based supervised learning techniques, decision tree regression, random forest, and gradient boosting, to systematically evaluate the online labor supply with 70 features related to climate, population, economics, education, health, language, and technology adoption. To provide machine learning explainability, SHAP, based on the Shapley values, was introduced to identify features with high marginal contributions. The top 5 contributing features indicate the tight integration of technology adoption, language, and human migration patterns with the online labor market supply.

Keywords

business, boosting, commerce and trade, digital divide, economics, ensemble learning, globalization, machine learning, random forest, social factors, statistical learning, sharing economy

1. Introduction

The gig economy is a system where people provide short-term goods and services. An example of the gig economy is freelance work, which can be completed independently or under a more prominent company acting as an intermediary platform. The online labor market is part of this gig economy, and it facilitates exchanges of all virtual services: software development, multimedia content, translation, and marketing support.

From 2016 to 2021, the online labor market has grown by over 10% annually (Stephany et al., 2021). This corroborates the idea that globalization increases the exchanges of products or services across countries (Friedman, 2005). The online labor market has gained greater visibility with the growth of internet accessibility and the COVID-19 pandemic from 2020, when there was a rise in the need for virtual services with limited face-to-face contact worldwide (Tan et al., 2021).

To better understand this growth, machine learning techniques were used to evaluate features associated with labor supply in the online labor market. Machine learning models were created by using data on the online labor market activities, climate, population, economics, education, health, language, and technology adoption over 5 years. With these data points, 6 models, multiple linear regression, Ridge, LASSO, decision tree regression, random forest, and gradient boosting, were trained, validated, and evaluated for factors related to the online labor market.

2. Data

Measurement of online labor market activity is based on the Online Labour Index collected from the Online Labour Observatory created by the International Labour Organisation and the Oxford Internet Institute (Stephany et al., 2021). This data on the online labor market supply were collected by examining the application programming interfaces from digital platforms or downloading the web user interface of 5 online labor platforms: Amazon Mechanical Turk, Upwork.com, Freelancer.com, Guru.com, and Peopleperhour.com. These were the top 5 platforms representing at least 70% of the total online labor platform traffic, according to Alexa.com, when the index was created in 2016. As of 2023, Freelancer.com is the only one of these 5 platforms that supports multiple languages, while the other 4 platforms only support

English. There are 992170 rows of daily online activities from June 16, 2017, to December 31, 2022 (Stephany et al., 2021).

The data gathered by the Online Labour Observatory was determined by looking at the number of open vacancies or projects for clients who wanted to hire workers. The occupation types of these vacancies were classified using machine learning based on specific keywords in the vacancy title or description. Employer country was estimated by taking a sample of vacancies from the two platforms, Guru.com and Upwork.com, that show country information, then weighing the samples to reflect all platforms' occupation distribution. The observations found in 203 countries were separated based on country, number of workers, number of projects, and occupation. Occupation was divided into six categories: clerical and data entry, creative and multimedia, professional services, sales and marketing support, software development and technology, and writing and translation. The online economy data were combined with data from the International Monetary Fund (International Monetary Fund, 2022), United Nations Development Programme (United Nations Development Programme, 2022), Oxford COVID-19 Government Response Tracker (Hale et al., 2021), United Nations Department of Economic and Social Affairs (United Nations Department of Economic and Social Affairs, 2023), World Factbook (Central Intelligence Agency, 2022), Ethnologue (Eberhard et al., 2019), Area Database of the Global Data Lab (Smits & Permanyer, 2019; Smits, 2016), World Bank (World Bank, 2020), and United Nations High Commissioner for Refugees (United Nations High Commissioner for Refugees, 2022) to create a dataset for modeling.

2.1. Features

Seventy features related to schooling, gross national income, labor force, technology subscriptions, tax rates, unemployment rates, inflation rates, migration, country expenditures, life expectancy, population, natural disasters, number of speakers of top languages, refugees, and COVID-19 were categorized for correlation analysis in this study ([Appendix A Table 1](#)).

Data were collected in Arabic, Chinese, English, French, Hindi, Russian, and Spanish since they were the official languages designated by the United Nations. The number of first language (L1) and second language (L2) speakers and the Expanded Graded Intergenerational Disruption Scale (EGIDS) (Lewis & Simons, 2010) of these languages across all countries was collected from Ethnologue to test the impact of languages on the growth of the online labor market (Eberhard et

al., 2019). The categorical EGIDS data were converted to numerical value in this study, similar to an approach taken in modeling of language adoption in the digital space (Kornai, 2013). Data on COVID-19 and internet users were also collected to test the potential increase in online labor market supply with increased isolation and internet users. For some attributes, data were missing in the latter years of 2021 or 2022. In these cases, data were taken from the most recent year recorded and was inputted for the following missing years.

2.2. Correlation

The combined dataset was evaluated for correlation among features. For linear models, if multiple features have high correlations, the model may become inaccurate and generate a larger error of sum squares (Bühlmann et al., 2013). The existence of correlated features could also increase the error of tree-based models by reducing the effectiveness of features that balance the other features with heavy correlation (Tolosi & Lengauer, 2011). The dataset was analyzed using Pearson correlation and Spearman's ranking correlation to identify correlated features.

The Pearson correlation coefficient was used to evaluate collinearity between 2 features. A negative Pearson correlation coefficient indicates an inverse relationship between two features, while a positive Pearson coefficient represents a positive relationship between two features. This Pearson correlation is based on the assumption of continuous features following a linear relationship and assumes each observation has a pair of values without accounting for outliers. The formula for the Pearson correlation coefficient r is computed using the equation:

$$r = \frac{n(\sum x_i y_i) - (\sum x_i)(\sum y_i)}{\sqrt{[n\sum x_i^2 - (\sum x_i)^2][n\sum y_i^2 - (\sum y_i)^2]}} \quad (1)$$

where n is the number of pairs and x, y are individual points.

For strongly correlated features paired with coefficients with absolute values above 0.7, clusters of features were identified using hierarchical clustering. Based on hierarchical clustering with Ward's linkage, a dendrogram that displays the high collinearity based on the Pearson correlation coefficient was created (**Figure 1**). In each cluster, one feature was chosen to represent the cluster while the rest of the features were removed from the dataset in this feature selection process.

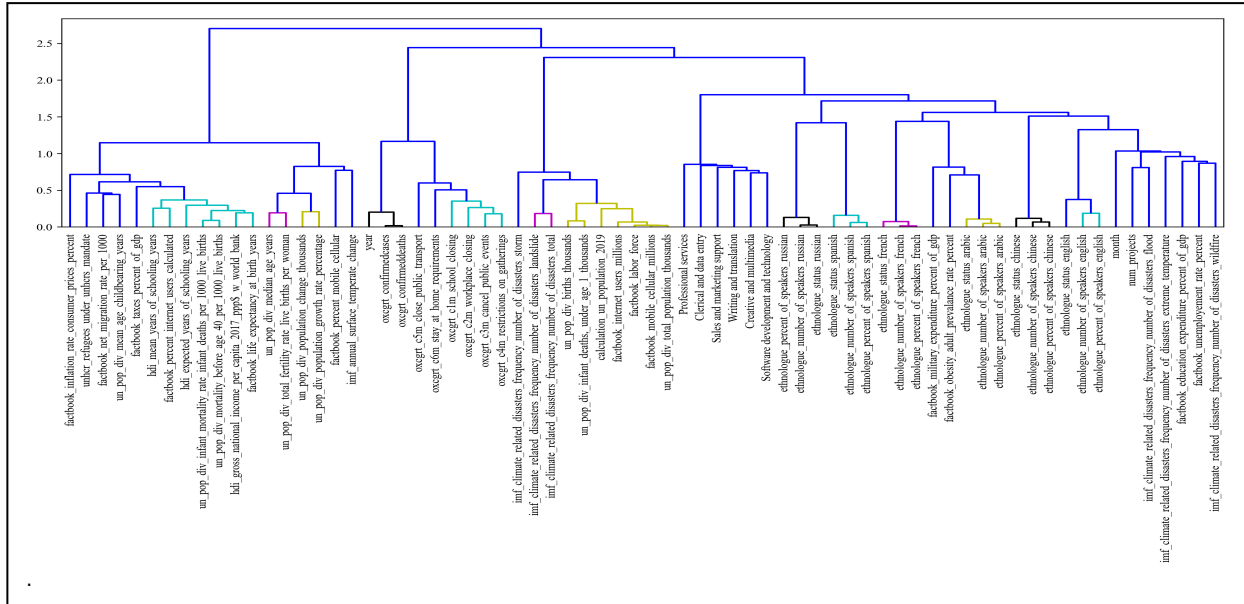


Figure 2. Dendrogram based on clustering with Spearman coefficient.

3. Models

Six models were used to identify the best match for the data. For all models, data were split into a training set and a testing set. The train set of data is used to create the model and train it to be more accurate, while the test set validates the model. A third of randomized data were used for testing the data, while the remaining two-thirds of data were used for training.

3.1. Accuracy Indicators

Accuracy of the models is measured by Mean Absolute Error (MAE), Root Mean Square Error (RMSE), and the coefficient of determination (R^2).

3.2. Linear Models (Multiple Linear Regression, Ridge Regression, LASSO Regression)

In this study, 3 linear models were used: multiple linear regression, ridge regression, and LASSO regression (James et al., 2023). The shrinkage penalties in the ridge and LASSO regressions are adjusted to tune the models.

3.3. Tree-based Models (Decision Tree, Random Forest, Gradient Boosting)

Three tree-based models, decision tree regression, random forest, and gradient boosting, were tested in this study (Breiman, 2001). To tune these models, maximum tree depth, number of

trees, and number of iterations were used on decision tree regression, random forest, and gradient boosting models respectively.

3.4. SHAP Values

SHAP, an abbreviation of SHapley Added exPlanations (Lundberg et al., 2020), values measure the feature importance of models based on the concept of Shapley values in game theory. Shapley values are calculated based on the marginal contributions of each component. For each component, a series of iterations was conducted by incrementally adding features to the model besides the component itself. The Shapley value is calculated by averaging the differences among all iterations without the component and the full model result. When features are heavily correlated, SHAP values stray further from their true Shapley value (Molnar, 2020).

4. Empirical Findings

In this study, all models were created using 3 datasets: the full dataset, the dataset with clusters based on Pearson coefficient analysis, and the dataset with clusters based on Spearman coefficient analysis.

4.1. Correlation Analysis with Pearson Coefficient

Correlation among features can create inaccurate modeling results in linear models. By analyzing the dataset based on the Pearson coefficient and creating a dendrogram based on Ward's linkage, 13 clusters of features were identified. Eleven out of the 13 new clusters had features formed within feature categories. Two clusters were identified to have features from multiple categories. One new cluster was identified based on the gross national income per capita, mean years of schooling, and expected year of schooling, while another cluster was formed based on the linear correlations based on mobile subscriptions, internet users, labor force, and population statistics. Although the correlation between years of schooling and gross national income per capita was found in this empirical analysis, knowledge capital is a more effective measure of education on economic output (Hanushek & Woessmann, 2020). In the scope of this study, gross national income per capita was selected during the feature selection process with the filtering method. A new dataset was created based on this analysis, and the number of features used has been reduced from 70 features to 48 features.

4.2. Correlation Analysis with Spearman Coefficient

Spearman coefficient analysis identified monotonic relationships among features. In the analysis, 13 clusters were identified. Two out of 13 clusters had features that spanned across multiple categories. One cluster was formed based on mean years of schooling, percent of internet users in the country, expected years of schooling years, infant mortality rate rates and gross national income, and life expectancy, while another cluster was formed based on births, infant deaths, population, internet users and mobile subscriptions, and number of people in the labor force. With hierarchical clustering, the features used in this new dataset were reduced from 70 to 39 features.

4.3. Evaluation of Linear Models

In **Table 2**, which displays each model's training and testing accuracy indicators, linear models had a higher error than the tree-based models. All linear models had low accuracy, and none of the linear models reached convergence after 1000 iterations during the training. Even though there may be a linear relationship between online labor market activity and features on an individual basis, online labor market activity does not have a linear relationship with features as an aggregate in this study.

Table 2. Comparison of MAE, RMSE, R^2 for all models.

Models	Model Accuracy Indicator					
	MAE Train	MAE Test	RMSE Train	RMSE Test	R^2 - Train	R^2 - Test
Multiple Linear Regression	877.750	875.346	2694.1257	2634.0875	0.3355	0.3338
Ridge Regression	877.726	875.322	2694.1257	2634.0874	0.3355	0.3339
LASSO Regression	877.180	874.844	2694.1662	2634.1458	0.3355	0.3338

Models	Model Accuracy Indicator					
	MAE Train	MAE Test	RMSE Train	RMSE Test	R ² - Train	R ² - Test
Decision Tree Regression	1.605	12.114	8.1017	177.5470	1.0000	0.9970
Random Forest	3.927	9.730	68.1244	150.8889	0.9996	0.9978
Gradient Boosting	877.750	875.346	2694.1257	2634.0875	0.3355	0.3338

4.4. Evaluation of Tree-Based Models

All 3 tree-based models, decision tree regression, random forest, and gradient boosting, were tuned by testing for the highest R² with combinations of model-specific parameters. The decision tree regression model's accuracy improved with increasing R² value and decreasing RMSE value when the maximum number of tree depths increased (Figure 3). Its value reached its highest when the depth was 21. For random forest models, the R² improved with the number of trees added to the model (Figure 4). R² was the highest when the number of trees was 34, and the maximum tree depth was 55.

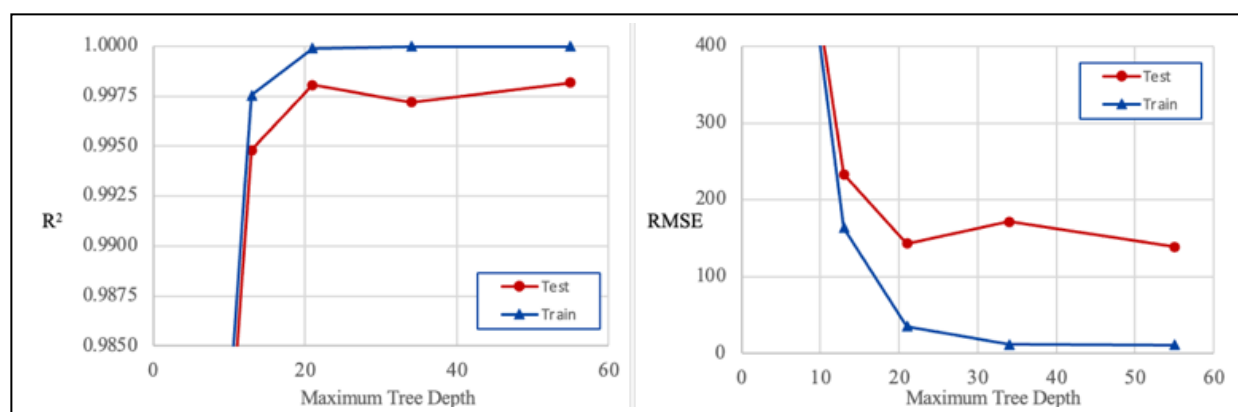


Figure 3. R² and RMSE as a function of maximum tree depth in decision tree regression model with feature selection using Spearman ranking coefficient.

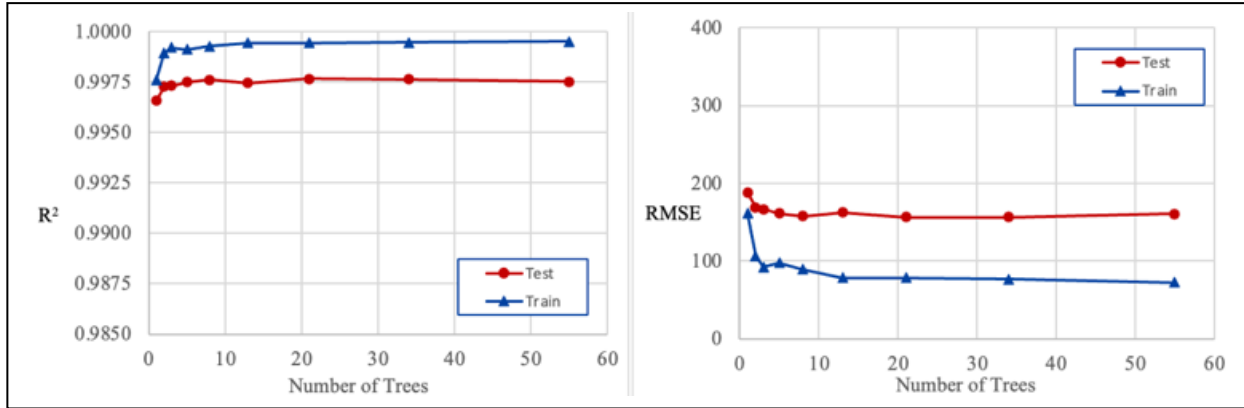


Figure 4. R^2 and RMSE as a function of number of trees in random forest model with feature selection using Spearman ranking coefficient.

In the case of gradient boosting modeling, improvements were found in R^2 with the number of iterative corrections being made to the model (**Figure 5**). Despite the iterative improvements, the RMSE of the gradient boosting models was still the highest among all tree-based models.

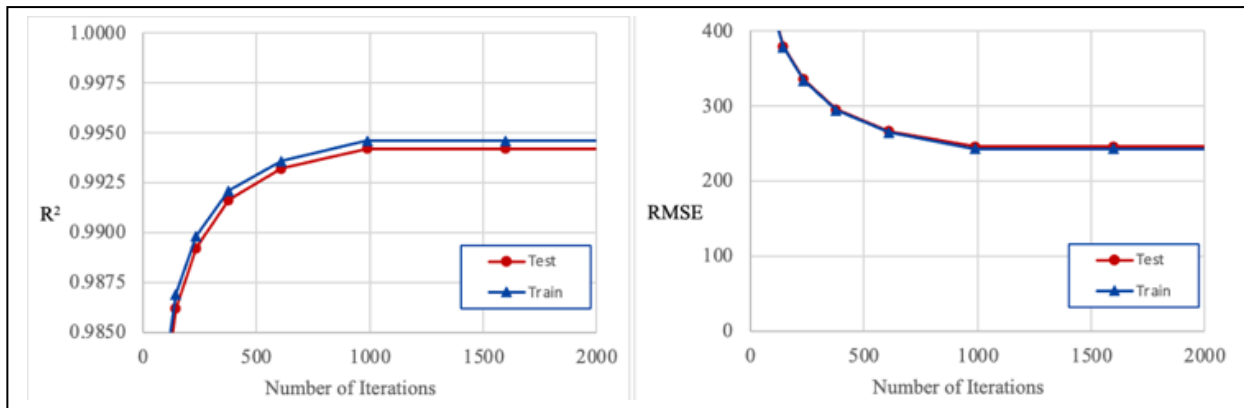


Figure 5. R^2 and RMSE as a function of number of iterations in gradient boosting model with feature selection using Spearman ranking coefficient.

Using these 3 parameters to establish the optimized models, tree-based models had higher accuracy and lower errors than the linear models (**Table 2**). The average testing MAE values for all tree-based models were 96.38% smaller than the testing MAE values for linear models. The average testing RMSE values for tree-based models was 93.21% smaller than the average testing RMSE values for linear models. For model fit, the average testing R^2 values for tree-based

models were 242.15% larger than that of linear models. The random forest model had higher accuracy than the LASSO model, as shown in the residual plot comparison (**Figure 6**).

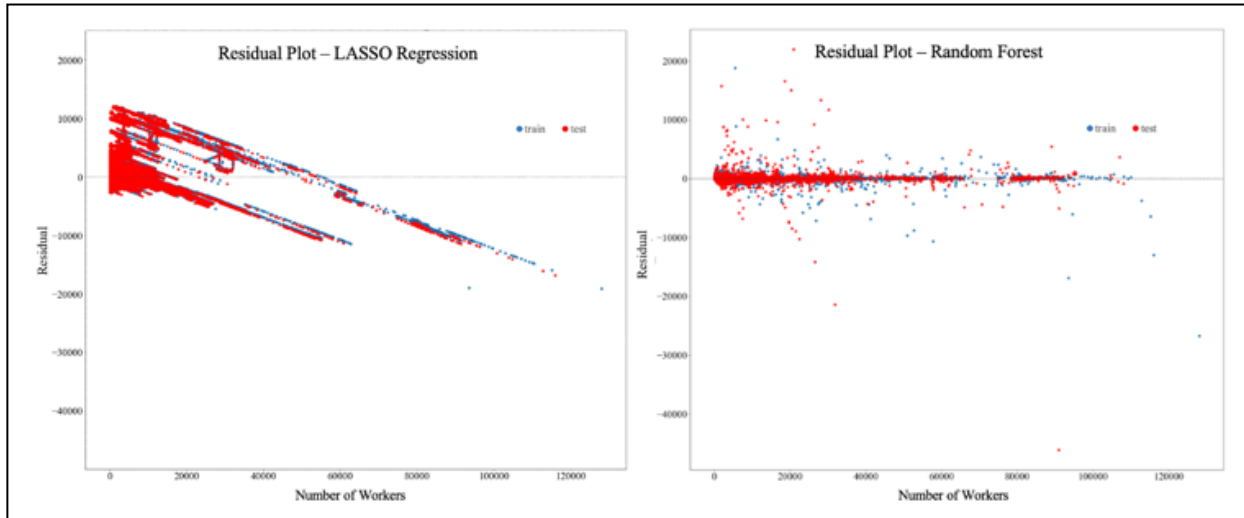


Figure 6. Residual comparison between LASSO and random forest models.

For tree-based models, highly correlated features co-existing in the dataset can affect the evaluation of the feature importance (Molnar, 2020). Feature selection based on Pearson and Spearman coefficients did not affect the R^2 of the tree-based models (**Table 2, Table 3, Table 4**). In comparison, for all 3 tree-based models, R^2 values of training data among all 3 datasets varied by less than 0.25%. Random forest models with the feature selections had the least variation in RMSE and MAE test values from the model with the original dataset. By applying feature selection, RMSE and MAE values of the random forest models varied by less than 5%.

Table 3. Model Accuracy: Dataset with feature selection based on Pearson coefficient ≥ 0.7

Models	Model Accuracy Indicator					
	MAE	MAE Test	RMSE	RMSE	R ² - Train	R ² - Test
	Train	Test	Train	Test		
Multiple Linear Regression	927.830	925.250	2772.7066	2706.4776	0.2962	0.2967
Ridge Regression	927.633	925.622	2771.4472	2709.0626	0.2967	0.2956

Models	Model Accuracy Indicator					
	MAE Train	MAE Test	RMSE Train	RMSE Test	R ² - Train	R ² - Test
LASSO Regression	927.625	925.614	2771.4471	2709.0619	0.2967	0.2956
Decision Tree Regression	0.990	11.048	6.9186	156.0806	1.0000	0.9977
Random Forest	3.873	9.494	68.9942	148.4736	0.9996	0.9979
Gradient Boosting	98.849	101.128	285.3918	271.5192	0.9925	0.9929

Table 4. Model Accuracy: Dataset with feature selection based on Spearman coefficient ≥ 0.7

Models	Model Accuracy Indicator					
	MAE Train	MAE Test	RMSE Train	RMSE Test	R ² - Train	R ² - Test
Multiple Linear Regression	1021.516	1020.021	2871.6414	2807.3630	0.2450	0.2433
Ridge Regression	1021.923	1018.888	2872.2688	2805.9725	0.2446	0.2443
LASSO Regression	1021.970	1018.936	2872.2688	2805.9730	0.2446	0.2443
Decision Tree Regression	17.027	24.707	35.4358	142.7089	0.9999	0.9980
Random Forest	4.04	10.165	72.7256	155.0682	0.9995	0.9977
Gradient Boosting	54.808	57.719	226.5599	224.8613	0.9953	0.9951

Comparing the accuracy indicators between testing and training data can identify the level of underfitting or overfitting among all the models. Among the 3 tree-based models, decision tree regression models have overfitting with large differences in RMSE and MAE values between test and train data and low RMSE and MAE training scores.

Gradient boosting models have the least overfitting among all tree-based models. It was the only tree-based model where the RMSE of the training data was higher than that of the testing data. While gradient boosting has the least overfitting among all tree-based models, it has the highest MAE and RMSE values and the lowest R^2 values.

Overall, among the 3 tree-based models, the random forest model was the most optimal model, which had lower error than the gradient boosting model and less overfitting than the decision tree regression model. Random forest model also had the least variation in R^2 , MAE, and RMSE after the model was simplified with correlation analysis.

4.5. SHAP Values

Ensemble models, like random forest and gradient boosting, have had challenges in their explainability. In the past despite their high accuracy (Peet et al., 2022). SHAP was used to address the challenges of explainability of these ensemble models. SHAP values of features in random forest model with feature selection based on Spearman correlation coefficients was computed (**Figure 7**). In order of importance, the top five features identified in the random forest model were the number of mobile subscriptions, the number of English speakers, the software development occupation category, the creative and multimedia occupation category, and the net migration rate.

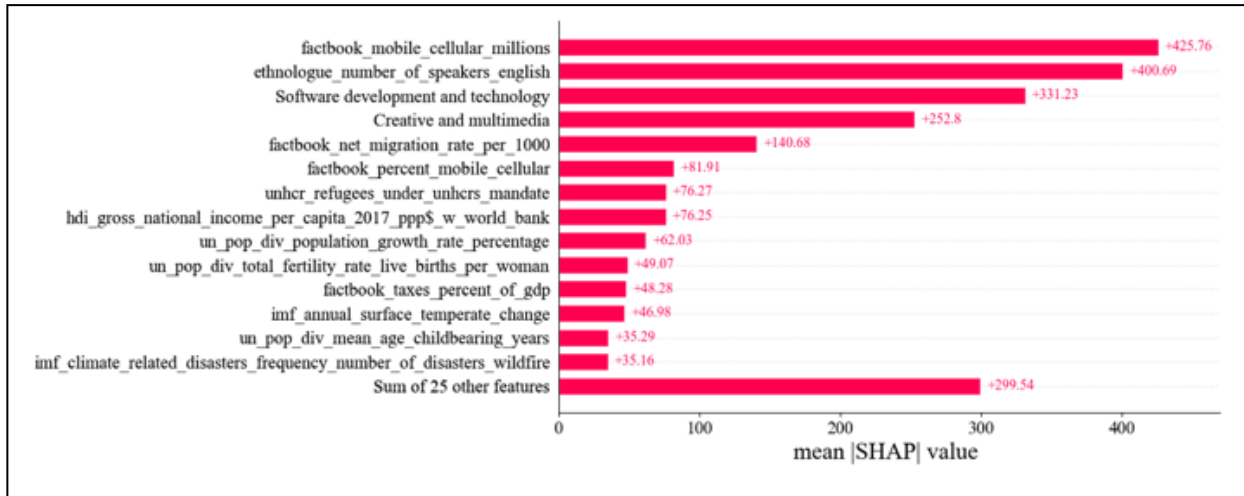


Figure 7. SHAP Values of random forest model with feature selection using Spearman coefficient.

5. Discussion

In the tree-based models with feature selection using Spearman correlation, none of the features was identified as the single most important contributing factor with a majority of all SHAP values. Although the number of English speakers, cellular mobile subscriptions, and software development were the top 3 features related to the online labor market supply, these factors represent 49% of the marginal contribution to the random forest model.

5.1. Mobile Cellular Subscriptions

When considering Pearson and Spearman correlation coefficients, clusters were formed based on the number of mobile subscriptions, internet users, labor force, and the total population, representing the strong adoption of the technology in everyday life where the mobile subscriptions rise in conjunction with the total population and internet users. The increased adoption of new technologies has a significant role in the growth of the online labor market. The supply of online labor activities in the top 10 countries, which represents 49.09% of all activities, has grown 26.56% during the period of this study from June 2017 to December 2022. These top 10 countries had 35.64% of all mobile subscriptions worldwide 2022. The result of this study reinforces the importance of information and communications technology (ICT) in creating access to economic opportunities for workers (Atasoy, 2013).

5.2. English Speakers

The online labor market spans across borders with participants from a large range of countries and territories. A common language is required for trade in this economy. Although English is not necessarily the required language for international trade, the ability for parties to communicate in a common language is more critical than trust and ethnicity (Melitz & Toubal, 2014). English's importance in all 3 tree-based models is a result of the need for a common language used in the demand and supply sides of this online labor market. From the demand side, the top 7 buying countries, representing 65.87% of all demand, have English as a national language. At the same time, from the supply perspective, the top 10 supplying countries, representing 49.09% of the 2022 online labor supply, had 61.03% of all first language (L1) and second language (L2) English speakers, as defined by Ethnologue (Stephany et al., 2021), (Eberhard et al., 2019). Although the number of English speakers is one of the top features based on SHAP values, it only represents 17% of the total contribution in the random forest model.

5.3. Software Development Category and Creative and Multimedia Category

The importance of the creative and multimedia occupation category is a result of an existing adoption of freelancing in the profession. Artists, who work in this occupation category, had been practicing freelancing for more than 20 years. From 2003 to 2015, in the United States, more than 30% of all artists were self-employed (Woronkiewicz & Noonan, 2019). Eighteen percent of all graphic designers, part of the artists group in this creative and multimedia occupation category, are self-employed workers as of 2023 in the United States (Bureau of Labor Statistics, U.S. Department of Labor, 2023). The online labor platforms and adoption of cloud applications supported freelancing in these creative and multimedia projects. (Sutherland & Jarrahi, 2017). Although the COVID-19 pandemic negatively affected the creative industry (Khlystova et al., 2022), this category remained the 2nd highest in demand in the online labor market (Stephany et al., 2021).

The importance of software development feature is a reflection of the demand for specific technical skills in the global market. Since 2016, software developer has been highlighted as an occupation with high demand across multiple countries (World Economic Forum, 2016, 2018, 2020, 2023). In 2022, software developer is one of the occupations with the most common labor shortage in Europe (World Economic Forum, 2023). While the online labor market provides

online opportunities to workers worldwide, the skills required for these opportunities need to be widely taught. National programs, such as the National Freelance Training Program by the Ministry of Information Technology & Telecommunications in Pakistan, have been teaching digital skills, such as software development and creative media, to meet the demand provided by freelance online gig work (Ministry of Information Technology & Telecommunication, 2023).

5.4. Migration Rate

In all 3 tree-based models, migration and the number of refugees had high SHAP values and were important for creating the models on the online labor market. Many migrants suffer in adapting to new countries where barriers limit job opportunities and render them uncertain (Altenried, 2021). The online labor market's low barriers of entry provide job opportunities to those in uncertain geographical positions. The online labor market often relies on migrants and their lack of alternative job options (Hackl & International Labor Organization, 2021). While the online labor market provides jobs, the income is often insufficient for a standalone job (Vernon et al., 2016). Although the number of refugees was one of the top 9 most contributing features, the lack of affordable mobile cellular access for refugees is a barrier to connecting with online economic activity (Vernon et al., 2016).

5. Conclusion

Random forest machine learning model is the optimal model to assess the factors that correlate with the labor supply in the online labor market. From the study, technology adoption, linguistics, and social factors are uncovered to be important to the overall supply at a global level.

Appendix

Table 1. List of all feature names used to create the models and their source.

Category	Features	
	Feature Name	Description
Climate	imf_annual_surface_temperate_change ¹	Temperature change with respect to a baseline climatology
Climate	imf_climate_related_disasters_frequency_number_of_disasters_extreme_temperature ¹	Number of days with extreme temperature during that year
Climate	imf_climate_related_disasters_frequency_number_of_disasters_flood ¹	Number of floods during that year
Climate	imf_climate_related_disasters_frequency_number_of_disasters_landslide ¹	Number of landslides during that year
Climate	imf_climate_related_disasters_frequency_number_of_disasters_storm ¹	Number of storms during that year
Climate	imf_climate_related_disasters_frequency_number_of_disasters_total ¹	Total number of disasters during that year
Climate	imf_climate_related_disasters_frequency_number_of_disasters_wildfire ¹	Number of wildfires during that year
Population	calculation_un_population_2019 ⁴	Population in the year 2019
Population	factbook_net_migration_rate_per_1000 ⁵	Net migration rate compares the difference between the number of persons entering and leaving a country during the year per 1,000 persons
Population	un_pop_div_births_thousands ⁴	Births (thousands)

Category	Features	
	Feature Name	Description
Population	un_pop_div_mean_age_childbearing_years ⁴	Mean age childbearing (years)
Population	un_pop_div_median_age_years ⁴	Median age (years)
Population	un_pop_div_population_change_thousands ⁴	Population change (thousands)
Population	un_pop_div_population_growth_rate_percentage ⁴	Population growth rate (percentage)
Population	un_pop_div_total_population_thousands ⁴	Total population (thousands)
Population	unhcr_refugees_under_unhcr's_mandate ⁹	number of refugees from country of origin
Labor	factbook_labor_force ⁵	Number of people in labor force
Economics	factbook_inflation_rate_consumer_prices_percent ⁵	Annual percent change in consumer prices with the previous year's consumer prices
Economics	factbook_taxes_percent_of_gdp ⁵	Total taxes and other revenues received by the national government, expressed as a percent of GDP
Economics	factbook_unemployment_rate_percent ⁵	Unemployment rate compares the percent of the labor force that is without jobs
Economics	hdi_gross_national_income_per_capita_2017_ppp\$w_world_bank ⁴	Gross National Income Per Capita using purchasing power parity rates in 2017
Education	factbook_education_expenditure_percent_of_gdp ⁵	Public expenditure on education as a percent of GDP
Education	hdi_expected_years_of_schooling_years ^{4,5}	Expected years of schooling (years)
Education	hdi_mean_years_of_schooling_years ^{4,7,8}	Mean years of schooling (years)

Category	Features	
	Feature Name	Description
Health	factbook_life_expectancy_at_birth_years ⁵	Life expectancy at birth compares the average number of years to be lived by a group of people born in the same year
Health	factbook_obesity_adult_prevalance_rate_percent ⁵	Adult prevalence rate gives the percentage of a country's population considered to be obese
Health	un_pop_div_infant_deaths_under_age_1_thousands ⁴	Infant Deaths, under age 1 (thousands)
Health	un_pop_div_infant_mortality_rate_infant_deaths_per_1000_live_births ⁴	Infant Mortality Rate (infant deaths per 1,000 live births)
Health	un_pop_div_mortality_before_age_40_per_1000_live_births ⁴	Mortality before Age 40, both sexes (deaths under age 40 per 1,000 live births)
Health	un_pop_div_total_fertility_rate_live_births_per_woman ⁴	Total Fertility Rate (live births per woman)
Health-Covid	oxcgrt_c1m_school_closing ³	Closing of schools for COVID-19
Health-Covid	oxcgrt_c2m_workplace_closing ³	Closing of workplaces for COVID-19
Health-Covid	oxcgrt_c3m_cancel_public_events ³	Cancellation of public events for COVID-19
Health-Covid	oxcgrt_c4m_restrictions_on_gatherings ³	Restrictions on gathering for COVID-19
Health-Covid	oxcgrt_c5m_close_public_transport ³	Closing of public transportation for COVID-19
Health-Covid	oxcgrt_c6m_stay_at_home_requirements ³	Stay at home requirements for COVID-19
Health-Covid	oxcgrt_confirmedcases ³	Cumulative number of reported COVID-19cases

Category	Features	
	Feature Name	Description
Health-Covid	oxcgrt_confirmeddeaths ³	Cumulative number of deaths attributed to COVID-19
Language	ethnologue_number_of_speakers_arabic ⁶	Number of Arabic language population
Language	ethnologue_number_of_speakers_chinese ⁶	Number of Chinese Mandarin language population
Language	ethnologue_number_of_speakers_english ⁶	Number of English language population
Language	ethnologue_number_of_speakers_french ⁶	Number of French language population
Language	ethnologue_number_of_speakers_russian ⁶	Number of Russian language population
Language	ethnologue_number_of_speakers_spanish ⁶	Number of Spanish language population
Language	ethnologue_percent_of_speakers_arabic ^{4,6}	% of Arabic speakers calculated from language population and total population in 2019
Language	ethnologue_percent_of_speakers_chinese ^{4,6}	% of Chinese Mandarin speakers calculated from language population and total population in 2019
Language	ethnologue_percent_of_speakers_english ^{4,6}	% of English speakers calculated from language population and total population in 2019
Language	ethnologue_percent_of_speakers_french ^{4,6}	% of French speakers calculated from language population and total population in 2019
Language	ethnologue_percent_of_speakers_russian ^{4,6}	% of Russian speakers calculated from language population and total population in 2019
Language	ethnologue_percent_of_speakers_spanish ^{4,6}	% of Spanish speakers calculated from language population and total population in 2019

Category	Features	
	Feature Name	Description
Language	ethnologue_status_arabic ⁶	EGIDS of Arabic in country
Language	ethnologue_status_chinese ⁶	EGIDS of Chinese in country
Language	ethnologue_status_english ⁶	EGIDS of English in country
Language	ethnologue_status_french ⁶	EGIDS of French in country
Language	ethnologue_status_russian ⁶	EGIDS of Russian in country
Language	ethnologue_status_spanish ⁶	EGIDS of Spanish in country
Military	factbook_military_expenditure_percent_of_gdp ⁵	Military expenditures as a percent of gross domestic product
Technology	factbook_internet_users_millions ⁵	Number of subscriptions within a country that access the Internet (millions)
Technology	factbook_mobile_cellular_millions ⁵	Number of mobile cellular telephone subscriptions (millions)
Technology	factbook_percent_internet_users_calculated ⁵	Number of internet subscriptions divided by population
Technology	factbook_percent_mobile_cellular ⁵	Number of mobile cellular subscriptions divided by population
Occupation	Writing and translation ¹⁰	Projects related to writing and translation
Occupation	Clerical and data entry ¹⁰	Projects related to clerical
Occupation	Creative and multimedia ¹⁰	Projects related to creative
Occupation	Software development and technology ¹⁰	Projects related to software development
Occupation	Professional services ¹⁰	Projects related to professional services
Occupation	Sales and marketing support ¹⁰	Projects related to sales and marketing support
Time	year ¹⁰	Year of online labor activity
Time	month ¹⁰	Month of online labor activity

Notes

¹ See International Monetary Fund (2022).

² See United Nations Development Program (2022).

³ See Hale et al. (2021).

⁴ See United Nations Department of Economic and Social Affairs (2023).

⁵ See Central Intelligence Agency (2023).

⁶ See Eberhard et al. (2019).

⁷ See Smits & Permanyer (2019); Smits (2016).

⁸ See World Bank (2020).

⁹ See United Nations High Commissioner for Refugees (2022).

¹⁰ See Stephany et al (2021).

References

- Altenried, M. (2021). Mobile workers, contingent labour: Migration, the gig economy and the multiplication of labour. *Environment and Planning A: Economy and Space*, 0(0).
<https://doi.org/10.1177/0308518x211054846>
- Atasoy, H. (2013). The Effects of Broadband Internet Expansion on Labor Market Outcomes. *ILR Review*, 66(2), 315-345. <https://doi.org/10.1177/001979391306600202>
- Breiman, L. (2001). Random Forests. *Machine Learning*, 45, 5–32.
<https://doi.org/10.1023/A:1010933404324>
- Bühlmann, P., Rütimann, P., van de Geer, S., & Zhang, C.-H. (2013). Correlated variables in regression: Clustering and sparse estimation. *Journal of Statistical Planning and Inference*, 143(11), 1835–1858. <https://doi.org/10.1016/j.jspi.2013.05.019>
- Bureau of Labor Statistics, U.S. Department of Labor. (2023). *National employment matrix 27-1024 Graphic designers* [Data set]. Bls.gov.
<https://data.bls.gov/projections/nationalMatrix?queryParams=27-1024&ioType=o>
- Central Intelligence Agency. (2022). *The world factbook*. Retrieved August 30, 2023, from <https://www.cia.gov/the-world-factbook/>.
- Eberhard, D., Simons, G., & Fennig, C. (2019). *Ethnologue: Languages of the World*. Ethnologue: Languages of the World, 22nd Edition.
- Friedman, T. L. (2005). *The world is flat: A brief history of the twenty-first century*. Macmillan.
- Hackl, A., & International Labor Organization. (2021). *Digital refugee livelihoods and decent work : towards inclusion in a fairer digital economy*. Retrieved August 30, 2023, from <https://policycommons.net/artifacts/1528335/digital-refugee-livelihoods-and-decent-work/2218020/>.

- Hale, T., Angrist, N., Goldszmidt, R., Kira, B., Petherick, A., Phillips, T., Webster, S., Cameron-Blake, E., Hallas, L., Majumdar, S., & Tatlow, H. (2021). A global panel database of pandemic policies (Oxford COVID-19 Government Response Tracker). *Nature Human Behaviour*, 5(4), 529–538. <https://doi.org/10.1038/s41562-021-01079-8>
- Hanushek, E. A., & Woessmann, L. (2020). Education, knowledge capital, and economic growth. *The Economics of Education*, 171-182. Academic Press.
<https://doi.org/10.1016/B978-0-12-815391-8.00014-8>
- International Monetary Fund. (2022). *Climate Change Indicators Dashboard* [Data set].
<https://climatedata.imf.org/pages/access-data>
- James, G., Witten, D., Hastie, T., Tibshirani, R., & Taylor, J. (2023). *An introduction to statistical learning: With applications in python*. Springer Cham.
- Khlystova, O., Kalyuzhnova, Y., & Belitski, M. (2022). The impact of the COVID-19 pandemic on the creative industries: A literature review and future research agenda. *Journal of Business Research*, 139, 1192–1210. <https://doi.org/10.1016/j.jbusres.2021.09.062>
- Kornai, A. (2013). Digital language death. *PloS One*, 8(10), e77056.
<https://doi.org/10.1371/journal.pone.0077056>
- Lewis, M. & Simons, G. (2010). Assessing Endangerment: Expanding Fishman’s Gids. *Revue roumaine de linguistique*, 55(2), 103–120. Retrieved August 30, 2023, from
<https://www.lingv.ro/RRL-2010.html>.
- Lundberg, S. M., Erion, G., Chen, H., DeGrave, A., Prutkin, J. M., Nair, B., Katz, R., Himmelfarb, J., Bansal, N., & Lee, S.-I. (2020). From local explanations to global understanding with explainable AI for trees. *Nature Machine Intelligence*, 2(1), 56–67.
<https://doi.org/10.1038/s42256-019-0138-9>

- Melitz, J., & Toubal, F. (2014). Native language, spoken language, translation and trade. *Journal of International Economics*, 93(2), 351–363.
<https://doi.org/10.1016/j.jinteco.2014.04.004>
- Ministry of Information Technology & Telecommunication. (2023, August 30). *NFTP | National Freelance Training Program*. <https://nftp.pitb.gov.pk>
- Molnar, C. (2022). *Interpretable Machine Learning: A Guide for Making Black Box Models Explainable* (2 ed.).
- Peet, E., Vegetabile, B., Cefalu, M., Pane, J., & Damberg, C. (2022). *Machine learning in public policy: the perils and the promise of interpretability*. <https://doi.org/10.7249/PEA828-1>
- Smits, J. (2016). GDL Area Database. Sub-national development indicators for research and policy making. *GDL Working Paper*, 16–101. Retrieved August 30, 2023, from <https://globaldatalab.org/asset/286/Smits%20GDL%20Working%20Paper%202016-101%20v360.pdf>.
- Smits, J., & Permanyer, I. (2019). The Subnational Human Development Database. *Scientific data*, 6, 190038. <https://doi.org/10.1038/sdata.2019.38>
- Stephany, F., Kässi, O., Rani, U., & Lehdonvirta, V. (2021). Online Labour Index 2020: New ways to measure the world's remote freelancing market. *Big Data & Society*, 8(2).
<https://doi.org/10.1177/205395172111043240>
- Sutherland, W., & Jarrahi, M. H. (2017). The gig economy and information infrastructure: The case of the digital nomad community. *Proceedings of the ACM on Human-Computer Interaction*, 1(CSCW), 1–24. <https://doi.org/10.1145/3134732>

- Tan, Z. M., Aggarwal, N., Cowls, J., Morley, J., Taddeo, M., & Floridi, L. (2021). The ethical debate about the gig economy: A review and critical analysis. *Technology in Society*, 65(101594), 101594. <https://doi.org/10.1016/j.techsoc.2021.101594>
- Tolosi, L., & Lengauer, T. (2011). Classification with correlated features: unreliability of feature ranking and solutions. *Bioinformatics*, 27(14), 1986–1994. <https://doi.org/10.1093/bioinformatics/btr300>
- United Nations Department of Economic and Social Affairs. (2022). *World Population Prospects 2022*. UN DESA Publications. Retrieved August 30, 2023, from <https://desapublications.un.org/publications/world-population-prospects-2022-summary-results>.
- United Nations Development Program. (2022). *Human Development Report 2021-22: Uncertain Times, Unsettled Lives: Shaping our Future in a Transforming World*. New York. Retrieved August 30, 2023, from <https://policycommons.net/artifacts/3533799/human-development-report-2021-22-human-development-reports/4335012/>.
- United Nations High Commissioner for Refugees. (2022). *Population figures* [Data set]. Unhcr.org. <https://www.unhcr.org/refugee-statistics/download/?url=HF39gS>
- Vernon, A., Deriche, K., & Eisenhauer, S. (2016). *Connecting refugees—how Internet and mobile connectivity can improve refugee well-being and transform humanitarian action*. Retrieved August 30, 2023, from <https://www.unhcr.org/media/connecting-refugees>.
- World Bank. (2020). *Education statistics - all indicators* [Data set]. https://databank.worldbank.org/indicator/SE.SCH.LIFE?id=c755d342&report_name=EdStats_Indicators_Report&populartype=series
- World Economic Forum. (2016). *The Future of Jobs 2016*. Retrieved August 30, 2023, from <https://www.weforum.org/reports/the-future-of-jobs-2016>.

World Economic Forum. (2018). *The Future of Jobs Report 2018*. Retrieved August 30, 2023, from <https://www.weforum.org/reports/the-future-of-jobs-report-2018>.

World Economic Forum. (2020). *The Future of Jobs Report 2020*. Retrieved August 30, 2023, from <https://www.weforum.org/reports/the-future-of-jobs-report-2020>.

World Economic Forum. (2023). *The Future of Jobs Report 2023*. Retrieved August 30, 2023, from <https://www.weforum.org/reports/the-future-of-jobs-report-2023>.

Woronkowicz, J., & Noonan, D. S. (2019). Who goes freelance? The determinants of self-employment for artists. *Entrepreneurship Theory and Practice*, 43(4), 651–672. <https://doi.org/10.1177/1042258717728067>