# EXPOSYS DATA LABS PROJECTS

**Title of The Report:**

profit Prediction the use of Regression fashions

**Subtitle:**

A Comparative evaluation of various Regression techniques for Predicting company earnings

**prepared with the aid of:**

Faiza Talamarla

**association:**

**department:** CSE

**university:** IIIT RGUKT RK VALLEY

**Date of Submission:**

1306/2024

# Table of contents:

# Predicting company profit with the use of Regression models

## abstract

accurately predicting company income is essential for making well-knowledgeable financial choices and strategic plans. in this undertaking, we employ various models everyday forecast corporations' earnings based on their spending on R&D, administration, and advertising and marketing. Our most important purpose is daily create more than one algorithms, verify their performance using key metrics, and decide the most appropriate version for precise profit prediction.

We work with a dataset that includes financial facts from 50 businesses, together with R&D expenses management charges, marketing expenditures, and income figures. The records undergoes preprocessing steps like cleansing, dealing with lacking values, and encoding express variables. via Exploraday-to-dayry facts evaluation (EDA), we benefit insights inday-to-day the relationships among these variables.

several regression models are applied in this studies: Linear Regression, Polynomial Regression, Random forest Regression, support Vecdailyr Regression, decision Tree Regression, and Gradient Boosting Regression. each version's overall performance is evaluated using metrics like R^2 rating, suggest Absolute mistakes (MAE), imply Squared mistakes (MSE), and Root mean Squared error (RMSE).

A comparison of the fashions suggests that Linear Regression offers the maximum accurate results for this dataset. Its predictions are each dependable and specific.

The outcomes of this project illustrate the effectiveness of regression fashions in predicting agency income and underscore the significance of selecting a version primarily based on comprehensive assessment criteria. future endeavors may also involve exploring advanced techniques and larger datasets every day similarly improve prediction accuracy.

# Introduction

### historical past:

profit prediction stands as a essential factor in crafting business techniques, allowing companies day-to-day make nicely-knowledgeable economic choices and strategize for future enlargement. particular profit projections useful resource inside the green allocation of resources, optimization of expenses, and ordinary enhancement of profitability. In this era of facts-centric choice-making, harnessing gadget daily strategies for income prediction offers awesome blessings as compared every dayeveryday processes.

### targets:

1. Crafting Regression models: The intention is everyday assemble and examine various regression fashions capable of forecasting employer income day-to-day on R&D spending, management charges, and advertising and marketing prices.

2. Assessing model performance: The plan entails evaluating extraordinary regression algorithms' performances the usage of general assessment standards.

3. Pinpointing the most excellent model: The goal is every day become aware of the maximum unique version for income prediction, offering actionable insights for businesses everyday refine their economic strategies.

## Scope:

This undertaking facilities on applying regression techniques every day expect the profits of fifty agencies. The dataset encompasses crucial everyday like R&D spending, administration charges, advertising charges, and income margins. The models produced in this analysis are aimed toward comprehending the connections between these variables and income tendencies whilst providing dependable income forecasts.

### importance:

This venture's significance lies in its capability every day beautify enterprise selection-making processes considerably. by using correctly projecting income, companies can:

better distribute assets for maximizing returns on investments.

Spot possibilities for optimizing expenditure.

Make informed strategic decisions daily on statistics every day boost up increase and profitability.

# Data Overview

The dataset we are the usage of here, every day as the "50 Startups" dataset, offers us data on a bunch of fifty businesses. It tells us about how a lot they spend unique areas and how much cash they make. This statistics from Kaggle and is commonly used for making predictions in gadget every day know everyday.

**overall Entries:** 50

**wide variety of categories**: four (R&D Spend, management fee, advertising Spend, profit)

This collection of statistics is high-quality for making predictions because it has numbers for things like R&D spending management charges, advertising fees, and profits. each row suggests what a enterprise spent and how much cash they made, which allows whilst seeking to bet income primarily based on those information.

**details about the information:**

**R&D Spend:** This element tells us approximately the investment agencies make in studies and improvement. it is a glide of numbers it's day-to-day critical due to the fact R&D greenbacks are frequently linked every day creating new stuff and earning extra dough.

**management value**: here, we see numbers day-to-day walking the at the back of-the-scenes commercial enterprise stuff. The expenses tied day-to-day managing operations shed light on whether all that admin spending connects with making greater financial institution.

**marketing Spend**: examine the cash splashed on advertising through each employer. money placed in daily advertising and marketing may have a massive impact on income and profits through attracting consumers and boosting revenues. it's another bunch of numbers every day play with.

**profit**: the important thing object in this dataset is earnings because it reveals how a whole lot moolah every enterprise takes home. Prediction models will goal everyday guess income relying on those 4 important capabilities.

**statistics exam and Key Stats:**

To apprehend facts, we calculate primary stats like average, center cost, spread, and scope for each class. those figures are awesome beneficial in spotting any weird or manner-out-there facts points that would mess with how accurate our version is making an attempt every day be.

**precis facts:**

- **R&D Spend:**

  - mean: 73721.62

  - Median: 73051.08

  - Deviation: 47483.06

  - range: 165349.2 - 0

- **administration cost:**

  - imply: 121344.64

  - Median: 122782.75

  - Deviation: 28017.eighty

  - variety: 182645.56 - 51283.14

- **marketing Spend:**

  - mean: 211025.09

  - Median: 212716.24

  - standard Deviation: 122290.31

  - variety: 471784.1 - 0

- **income:**

  - imply: 112012.64

  - Median: 107404.34

  - Deviation: 40306.18

  - range: 192261.83 - 14681.4

# data Preprocessing

**data cleaning:**

facts cleaning plays a vital position in preparing the dataset for gadget daily models. This process entails addressing lacking values, removing duplicates, and rectifying any inconsistencies or errors gift in the information.

to start with, let's search for any missing values within the dataset. it's lucky that the furnished dataset includes no lacking values, with every function having 50 non-null entries.

next, we want every day everyday for replica entries everyday maintain data integrity. Upon examination, no duplicate rows have been detected inside the dataset.

in terms of handling outliers, they have been diagnosed the use of field plots. apparently, these outliers have been retained within the dataset every day make certain facts integration remains intact.

transferring on to encoding specific variables inside the "50 Startups" dataset. fortuitously, there aren't any express variables requiring encoding as all features are numerical and continuous, making them suitable for regression evaluation.

thinking about that features within the dataset range in scales, it is crucial to apply scaling daily make certain all functions contribute similarly day-to-day the evaluation. Standardization (or Z-score normalization) is utilized .

this change results in facts with a median of 0 and a trendy deviation of 1. Such standardized information can beautify the overall performance of sure machine every day algorithms.

day-to-day investigate regression model performance efficaciously, the dataset is split indaily schooling and trying out units:

 **training Set**: Comprising 80% of the data

 **testing Set:** Comprising 20% of the records

This cut up ensures that the model is skilled on a sufficient quantity of records whilst also having unseen information for comparing its overall performance.

**Python code:**

from sklearn.model_selection import train_test_split

 define capabilities (X) and target (y)

X = df[['R&D Spend', 'Administration', 'Marketing Spend']]

y = df['Profit']

Splitting the dataset inevery day the schooling set and take a look at set

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=0)

# Exploratory Data Analysis(EDA):

Diving Deep in every day statistics: An Exploration with EDA

within the subject of records science, every assignment embarks on a adventure of exploration via facts evaluation (EDA). This vital segment serves as a basis, thoroughly assessing the fundamental traits of the facts. EDA makes use of a various set of visible tools daily animate the narrative hidden inside the statistics. It reveals underlying tendencies, identifies anomalies, and carefully tests theories – all pushed via the effective mix of precis information and compelling graphical illustrations.

in this enterprise, we delve  day an EDA journey with the "50 Startups" dataset. Our purpose: every day find the distribution of every characteristic and shed mild on its complex dating with the primary objective – maximizing revenue.

The impact of visible representation:

statistics visualization emerges as a pivotal daily, permitting us every day decode the tricky interaction amongst exceptional variables inside the dataset. It unveils concealed patterns, uncovers exciting tendencies, and brings outliers lurking in obscurity everyday light. by using leveraging visualization competencies, we go beyond easy precis records limitations, gaining a profound know-how that propels us ahead.

Hisday-to-daygrams: Revealing the function panorama

Our excursion commences with, important instruments for illustrating person feature distributions. those visual wonders provide insight day the information's critical inclinations, spread volume, and captivating shapes it embodies. Hisevery daygrams empower us to comprehend every feature's essence, laying a foundation for similarly investigation.

## Correlation evaluation:

Correlation evaluation is used every day renowned the strength and direction of relationships between numerical features. The correlation coefficient stages from -1 to one, wherein:

- "1" shows an excellent positive correlation

- "-1" shows a really perfect bad correlation

- "0" indicates no correlation

### Correlation Matrix

A correlation matrix gives a comprehensive view of the pairwise correlations among features. It enables in figuring out which features are most as it should be every dayeveryday the target variable (earnings).

### Python code:

```
correlation_matrix = df.corr()

sns.heatmap(correlation_matrix, annot=authentic, cmap='coolwarm')

plt.name('Correlation Matrix')

plt.show()
```

# Regression models :

Regression models intention everyday forecast non-stop target variables day-to-day on one or extra predict everyday variables. in this examine, numerous regression techniques are assessed every day discover the most reliable version for predicting employer profit based on elements everyday inclusive of R&D Spend, administration price, and advertising charges.

## Linear Regression:

explanation of Linear Regression model:

Linear regression serves as a statistical approach every day model the connection between a based variable (earnings) and unbiased variables (This model entails fitting a linear equation every day observed facts represented via:

$$y=\beta_0+\beta_1x_1+\beta_2x_2+\cdots+\beta_nx_n$$

• y denotes the predicted price (profit),

• $\beta_0$ represents the intercept,

• $\beta_1, \beta_2,\ldots, \beta_n$ imply coefficients for predict $x_1, x_2,\ldots, x_n$.

Alignment with data and Assumptions:

Linear regression accomplishes records alignment by using minimizing the sum of squared variations among real values and people anticipated with the aid of the version using the least squares technique. Key assumptions involve linearity in relationships among dependent and impartial variables, independence of observations from each different, steady errors variance (homoscedasticity), and typically dispensed residuals.

## Polynomial Regression:

Definition and alertness in Dataset Context:

Polynomial regression extends linear regression by modeling established and independent variables' dating as an nth-degree polynomial. through introducing polynomial phrases ineveryday the linear regression equation it captures non-linear relationships as represented by:

$$y=\beta_{zero}+\beta_1x+\beta_2x_2+\cdots+\beta_nx_n+\epsilon$$

## Random forest Regression:

### Conceptualization of Random forest Regression:

Random woodland regression adopts an ensemble every day approach employing a couple of selection timber for predicting target variables. each tree in the wooded area operates on a random subset of information and functions; afterward, predictions from all character bushes are averaged for very last prediction consequences. This approach mitigates overfitting issues enhancing model generalization abilities.

**benefits & Constraints:**

**advantages:**

* effectively manages tremendous datasets with better dimensions.

* provides insights on feature importance.

* Resilient in opposition to outliers & missing information.

Drawbacks:

* Computational demands could be excessive.

* performance may falter with very diminutive datasets.

## Support vector regression:

### overview & Implementation info of SVR:

aid  Regression falls below guide  device (SVM) umbrella intended for regression tasks optimizing hyperplane alignment within an mistakes margin whilst retaining version flatness. aimed at minimizing deviations among predicted values versus real figures inside a chosen margin (epsilon).

### Inclusion of Kernel type & Implications on version performance:

utilizing various kernel capabilities in SVR deal with non-linear relationships inside datasets efficaciously inclusive of Linear Kernel (for linearly separable data), Polynomial Kernel (for polynomial connections), Radial foundation feature (RBF) Kernel acceptable for complicated facts family members scrutiny. extensively favored in tasks daily its adaptability in taking pictures non-linear styles.

## decision Tree Regression :

### deciphering decision Tree Regression Operations:

using selection Tree Regression employs non-parametric modeling incorporating tree-like structures guiding predictive approaches reducing records subsets predicated on function

values for that reason constructing nodes & branches formed via certain criteria e.g., suggest squared error obtaining anticipated value outputs shaped at leaf nodes.

**Node Splitting techniques Embraced:**

regarding node splitting avenues in decision timber encompass measuring cumulative squared difference amid actual as opposed to projected figures via imply Squared blunders (MSE) or assessing average absolute discrepancy thru mean Absolute mistakes (MAE).

**Gradient Boosting Regression :**

**Elucidating Gradient Boosting Regression processes:**

Gradient Boosting mechanism harnesses ensemble day-to-day ethos sequencing more than one weak novices i.e., selection trees sequentially correcting earlier errors gathering advanced overall performance iteratively focusing on residuals refining precision richness bolstering robustness highlighting aggregate effort pooling weak learners' predictions harmoniously intensifying accuracy assurances.

# Model Evaluation

**performance Metrics:**

The performance of every model is evaluated the usage of the following metrics:

- **R² score**: This metric quantifies the distinctive feature of match of the model. It shows the proportion of the variance within the based variable this is predictable from the impartial variables.

- **mean Absolute errors (MAE):** Measures the common value of the errors in a set of predictions, without thinking about their course.

- **mean square error (MSE):** Measures the common of the squares of the errors—this is, the common squared difference among the expected values and the actual value.

- **Root mean square error (RMSE):** The rectangular root of the MSE, which measures the average magnitude of the error in the equal devices because the expected variable.

every version is educated on a schooling set and evaluated on a test set using five-fold cross-validation. This make certain that the fashions are trained and tested on different subsets of the records everyday keep away from overfitting.

 deciding on the best version

on this phase, we examine the performance of various regression fashions and identify the every dayryeveryday model for predicting employer income day-to-day on R&D Spend, management cost, and marketing Spend.

# Choosing the best model

**version performance assessment:**

The performance of every regression model is evaluated the use of the following metrics:

- R² score: Measures the proportion of the variance inside the structured variable (income) that is predictable from the unbiased variables (R&D Spend, management, advertising Spend). A better R² rating suggests a better healthy.

- mean Absolute blunders (MAE): Measures the average importance of the errors in a set of predictions, with out considering their direction.

- mean Squared error (MSE): Measures the average of the squares of the mistakes—that is, the common squared distinction among the envisioned values and the actual price.

- Root mean Squared error (RMSE): The rectangular root of the MSE, which measures the common magnitude of the mistake in the identical gadgets as the expected variable.

## summary of version performance

beneath is a precis table of the overall performance metrics for each regression model evaluated on the take a look at dataset:

| Model | R^2 | MAE | MSE | RMSE |
|---|---|---|---|---|
| Linear Regression | 0.900842 | 7499.964302 | 157193060.532141 | 12537.665673 |
| Polynomial Regression | 0.888484 | 7506.064597 | 90304882.2901 | 9502.888103 |
| Random Forest Regression | -0.725684 | 27019.127374 | 1397444892.075497 | 37382.414209 |
| Support Vector Regression | -0.180405 | 22851.85 | 955882952.97824 | 30917.356824 |
| Decision Tree Regression | -1.664267 | 43416.186 | 2157502424.3307 | 46448.922747 |
| Gradient Boosting Regression | -1.299111 | 37384.125923 | 1861802150.83522 | 43148.605433 |

**model choice:**

**Linear Regression** emerges as the daily performer among numerous fashions in forecasting company profits, showcasing a superior R² score of {{ linear_r2 }} and the maximum minimum RMSE of {{ linear_rmse }}. This underscores the efficacy of the linear version in hanging a harmonious stability between accuracy and simplicity inside the dataset.

In summary, through an intensive scrutiny of evaluation metrics and visible examinations, Linear Regression would be the every day desire for prognosticating organization earnings using the "50 Startups" dataset. This model gives a lucid interpretation of the correlation among R&D Spend, management price, advertising Spend, and profit, rendering it apt for guidance commercial enterprise decisions and strategic blueprints.

# Implementation

This segment delineates the execution of our system studying enterprise with the purpose of predicting enterprise profits hinged on R&D Spend, management price, and advertising Spend. The implementation encompasses statistics preprocessing, exploradailyry records analysis (EDA), version schooling throughout various regression frameworks, performance assessment, and outcomes visualization. beneath are the sequential steps along corresponding code snippets to demonstrate the holistic manner.

**equipment and Libraries hired:**

For this mission's execution, we leveraged the following Python libraries:

**pandas**: Facilitating statistics manipulation and evaluation.

**numpy:** permitting numerical computations.

**matplotlib and seaborn**: Orchestrating data visualization.

**scikit-examine**: Powering device every day knoweveryday algorithms and version appraisal

## information Loading and Inspection

We start by loading the dataset and inspecting its structure.

**Python code:**

Load the dataset:

information = pd.read_csv('/mnt/data/50_Startups.csv')

show the primary few rows of the dataset:

print(information.head())

show precis statistics:

print(facts.describe())

## information Preprocessing

Preprocessing the facts is essential every day ensure the model plays well. This involves handling missing values, encoding specific variables, and splitting the records ineveryday training and testing sets.

**Python code:**

coping with missing values (if any)

facts.fillna(data.mean(), inplace=true)

## Encoding express variables (if any)

on this dataset, there are not any express variables that want encoding.

**Python code:**

```
X = information[['R&D Spend', 'Administration', 'Marketing Spend']]

y = information['Profit']
```

**model evaluation**

The overall performance of each and each model turned into evaluated using the following metrics: suggest Absolute blunders (MAE), suggest Squared error (MSE), Root imply Squared mistakes (RMSE), and R-squared score (R^2).

**Python code:**

```
Predicting the consequences

y_pred_linear = linear_regressor.predict(X_test)

y_pred_poly = poly_regressor.expect(X_poly_test)

y_pred_rf = random_forest_regressor.predict(X_test)

y_pred_svr = svr_regressor.expect(X_test)

y_pred_dt = decision_tree_regressor.are expecting(X_test)

y_pred_gb = gb_regressor.expect(X_test)
```

**Calculating performance metrics**

```
fashions = {

    'Linear Regression': y_pred_linear,

    'Polynomial Regression': y_pred_poly,

    'Random wooded area Regression': y_pred_rf,

    'guide Vecevery dayr Regression': y_pred_svr,

    'selection Tree Regression': y_pred_dt,
```

```
    'Gradient Boosting Regression': y_pred_gb
}
for name, y_pred in models.gadgets():
    print(f"{name} performance:")
    print(f"MAE: {mean_absolute_error(y_test, y_pred)}")
    print(f"MSE: {mean_squared_error(y_test, y_pred)}")
    print(f"RMSE: {np.sqrt(mean_squared_error(y_test, y_pred))}")
    print(f"R^2 score: {r2_score(y_test, y_pred)}")
    print("-" * 30)
```

# Results and discussion

## summary of Findings:

The undertaking aimed everyday are expecting organization income the usage of R&D Spend, administration price, and advertising and marketing Spend thru diverse regression models. primary outcomes are:

1. **Linear Regression:**
   - MAE: 7499.964302
   - MSE: 157193060.532141
   - RMSE: 12537.665673
   - R^2: 0.900842

2. **Polynomial Regression:**
   - MAE: 7506.064597
   - MSE: 90304882.2901

- RMSE: 9502.888103

  - R^2: 0.888484


3. **Random woodland Regression:**

  - MAE: 27019.127374

  - MSE: 1397444892.075497

  - RMSE: 37382.414209

  - R^2: -0.725684


4. **support Vector Regression:**

  - MAE: 22851.eighty five

  - MSE: 955882952.97824

  - RMSE: 30917.356824

  - R^2: -zero.180405


5. **decision Tree Regression:**

  - MAE: 43416.186

  - MSE: 2157502424.3307

  - RMSE: 46448.922747

  - R^2: -1.664267


6. **Gradient Boosting Regression:**

  - MAE: 37384.125923

  - MSE: 1861802150.83522

  - RMSE: 43148.605433

  - R^2: -1.299111

**Linear Regression emerged as the first-rate model with an R^2 score of 0.90 and the lowest MAE of 7499.96**

**Interpretation of consequences:**

**Implications**

agencies can benefit from precise earnings predictions daily optimize spending in R&D, administration, and marketing. Strategic making plans every more effective with this version by means of pinpointing regions that have an impact on earnings, leading day-to-day higher funding choices. managing coins flows and improving productiveness day-to-day more doable while earnings prediction is accurate.

in the pursuit of challenge objectives, more than one regression fashions have been developed and evaluated. Linear Regression out as the most appropriate version aligning with the assignment's dreams. The consequences received from the assignment aid well-informed financial planning and strategic choice-making.

This mission the effectiveness of regression fashions in forecasting agency earnings based on R&D Spend, management cost, and advertising and marketing Spend. Linear Regression emerged as the superior model supplying precise predictions and interpretability for realistic commercial enterprise use.

Strategic monetary decision-making daily on information-pushed insights profits significance via these findings. with the aid of optimizing investments in R&D and advertising and marketing at the same time as ensuring efficient administrative spending, groups can raise profitability and acquire justifiable increase.

# Conclusion

## Key Findings:

1. model overall performance: Linear Regression changed into chosen as the best model for predicting organization income everyday its excessive R² score and favorable errors metrics after comparing multiple fashions.

2. version Interpretability: apart from accurate predictions, Linear Regression additionally presents interpretability making it an ideal choice for enterprise decisions.

3. feature significance: The observe revealed that R&D Spend has the most big effect on earnings accompanied by means of advertising Spend and management cost.

guidelines stemming from those findings suggest focusing greater on increasing R&D investments for better profitability, strategically allocating sources in advertising and marketing campaigns day-to-day pressure sales, and optimizing administrative expenses for stepped forward profitability.

**barriers & future paintings:**

1. Dataset size: With simplest 50 observations within the dataset, the generalizability of findings is restricted calling for a larger dataset in future work everyday beautify version robustness.

2. capabilities: Incorporating extra functions like industry kind, geographical location, and marketplace situations ought to offer deeper insights indaily earnings determinants.

3. superior techniques: Exploring superior machine learning strategies which includes neural networks and ensemble techniques holds promise to improve prediction accuracy.

# References

1. Dataset: "50 Startups" dataset from Kaggle day-to-day at [Kaggle Datasets](https://www.kaggle.com/farhanmd29/50-startups).

2. Géron, A. (2017). fingers-On machine every day with Scikit-analyze, Keras, and TensorFlow: standards, tools, and techniques every day build smart systems.

3. Raschka, S., & Mirjalili, V.. (2019). Python system daily: system every day and Deep mastering with Python, scikit-learn, and TensorFlow 2. Packt Publishing.

# Appendices:

**A - summary of facts functions:**

 R&D Spend: investment in studies & development.

administration cost: fees daily administrative obligations.

marketing Spend: Expenditure on advertising and marketing sports.

earnings: The ensuing take advantage of those prices.

## B - Preprocessing summary:

records cleansing concerned disposing of duplicates & irrelevant entries.

lacking values have been crammed the use of column way.

categorical variables had been encoded the usage of one-warm encoding.

## C - model performance Metrics precis:

A evaluation of overall performance metrics for each regression model is offered here.

| Model | R^2 | MAE | MSE | RMSE |
|---|---|---|---|---|
| Linear Regression | 0.900842 | 7499.964302 | 157193060.532141 | 12537.665673 |
| Polynomial Regression | 0.888484 | 7506.064597 | 90304882.2901 | 9502.888103 |
| Random Forest Regression | -0.725684 | 27019.127374 | 1397444892.075497 | 37382.414209 |
| Support Vector Regression | -0.180405 | 22851.85 | 955882952.97824 | 30917.356824 |
| Decision Tree Regression | -1.664267 | 43416.186 | 2157502424.3307 | 46448.922747 |
| Gradient Boosting Regression | -1.299111 | 37384.125923 | 1861802150.83522 | 43148.605433 |

## D: tools and Libraries

- Programming Language: Python

- Libraries: pandas, numpy, matplotlib, seaborn, scikit-analyze

## E: extra Notes

- Dataset supply: The dataset changed into amassed from [insert supply if now not daily in the foremost

- surroundings: The evaluation become carried out in a Jupyter pocket book environment.