

A Guide to Elevated Global Happiness

Nicole Huang and Aiden Kim

True or False

“Happiness is directly related to economic growth/wealth”

Self-reported life satisfaction vs. GDP per capita, 2024

Self-reported life satisfaction is measured on a scale¹ ranging from 0-10, where 10 is the highest possible life satisfaction. GDP per capita is adjusted for inflation and differences in living costs between countries.

Life satisfaction (0-10)

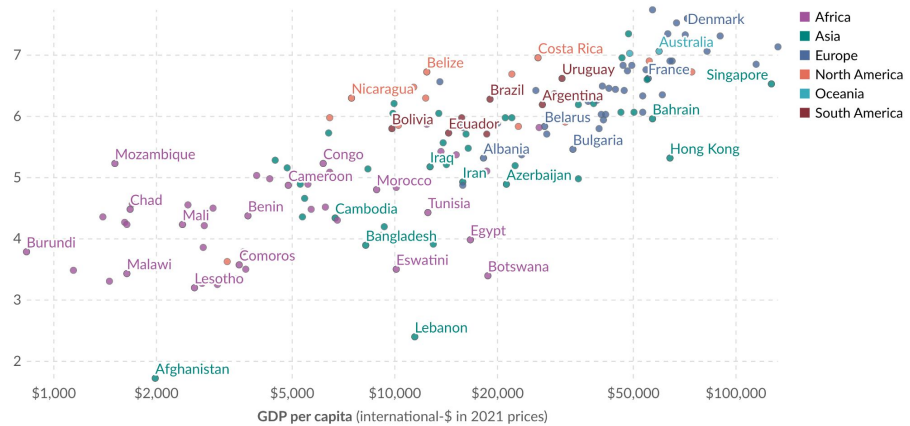


Figure created by [OurWorldinData.org](https://ourworldindata.org) based on data collected by Wellbeing Research Center (2025), compiled by the World Bank (2025)

Our Motivation

- Is it really the most definitive feature for happiness levels?
- Happiness is subjective (hard to measure and define)
- Global happiness provides insight to create informed policy decisions that cater to enhance people's lives

A decorative border surrounds the central text box, featuring various colored smiley faces (pink, yellow, green, blue) and small solid dots of the same colors. The background is a light gray grid.

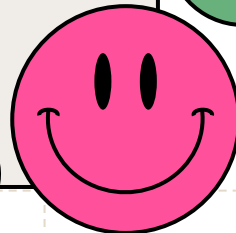
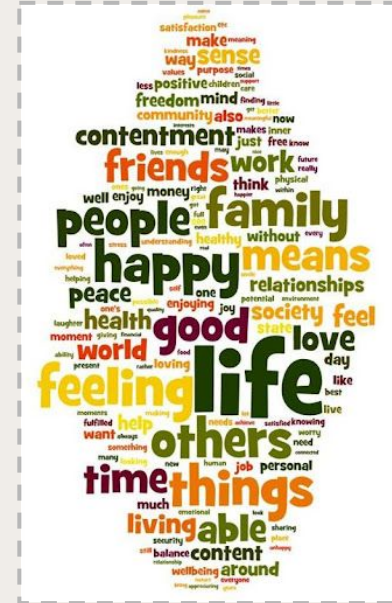
Scientific Question:

What feature is most directly predictive
of high happiness levels?

We Predicted perceived freedom closely associated with higher
global happiness levels – live their lives however they want.

Table of Contents

- 01 Data
- 02 Methods
- 03 Final results and analysis
- 04 Conclusion





01 Data

Dataset: “World Happiness Report”

Dataset of reported global happiness score with 21 associated features by various countries between 2005-2024.

- Source: kaggle.com
- Size: 4000x24 (n=4000 p=24)
- **Training data:** selected dataset from years 2005 (earliest)



Dataset included the following features:

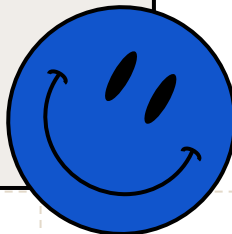
y-label

Country, Year, Happiness_Score, GDP_per_Capita, Social_Support, Healthy_Life_Expectancy, Freedom, Generosity, Corruption_Perception, Unemployment_Rate, Education_Index, Population, Urbanization_Rate, Life_Satisfaction, Public_Trust, Mental_Health_Index, Income_Inequality, Public_Health_Expenditure, Climate_Index, Work_Life_Balance, Internet_Access, Crime_Rate, Political_Stability, Employment_Rate

Features

Important Feature Data Types

- **Happiness_Score:** float between the scale of 1-10 (**label**)
- **Statistical Data** ie GDP_per_Capita, Health_Life_Expectancy, Population
- **Scaled ratings** (between 0-1) ie Social Support, Freedom, Public Trust



Data Cleaning

01

Removing Data

Removed all datasets that had null values.
(not much change)

02

Binary conversion

- Wanted binary labels
(low/high happiness levels)
- Threshold: 0.5

03

Features excluded

- Country
- Year
- Happiness_Score
(label)

Finalized dataset sorted into dictionary organized by year.

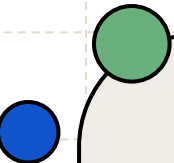
2005 dataset used as training data.



02

Methodology

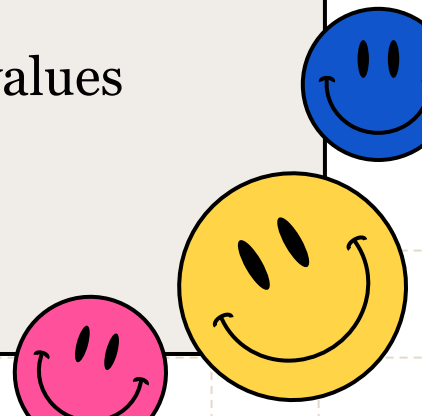
+ Initial Results



Our goal Identify the **best feature** that can closely **predict** high happiness levels.

Method: We want to find
lowest entropy -> highest gain = best feature

Then what? Compare predicted and actual gain values
to determine predicted accuracy.



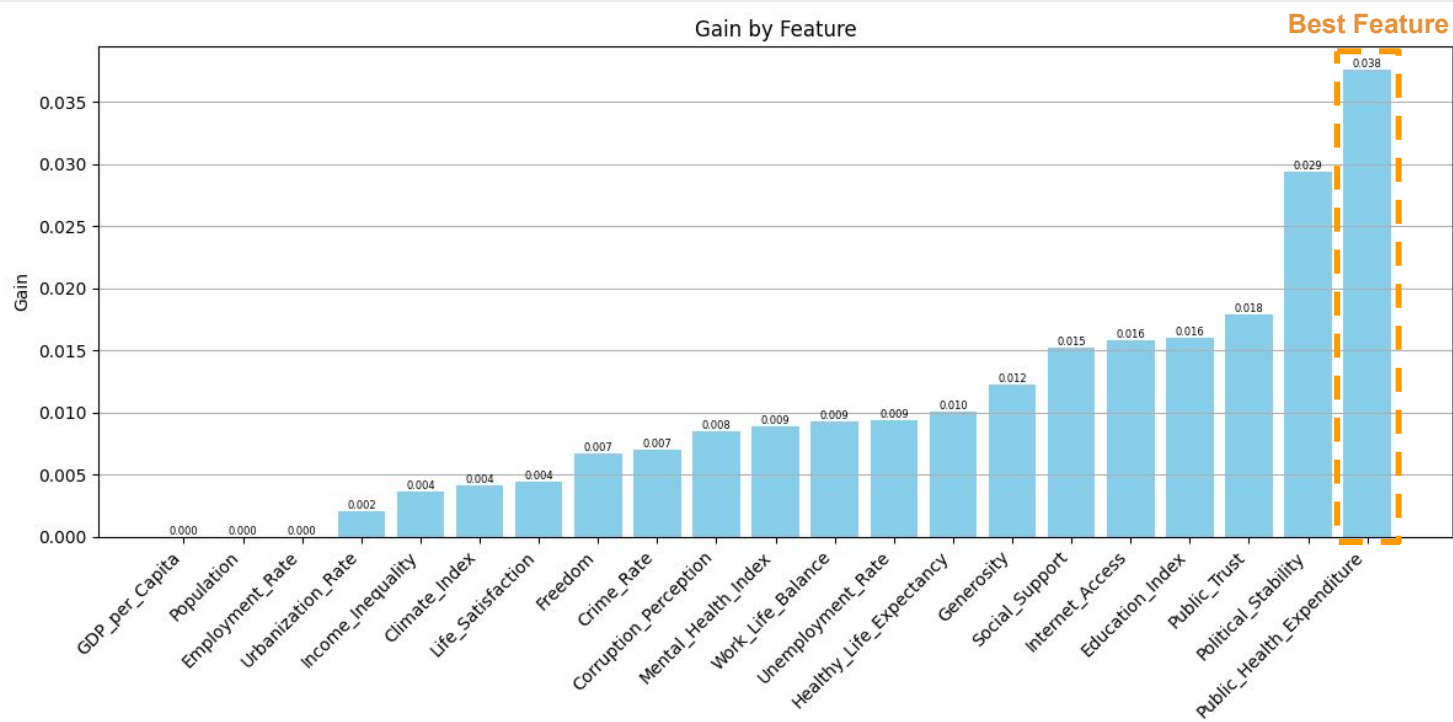
“Training Process”

Gain = total entropy - feature entropy

- For each feature, we calculated the gain value.

Highest gain value = best feature

- Store the **best feature** → **pred_feat_entropy**



Data Collection Process

- Assume best_feature consistent throughout all datasets
- Calculate gain of best feature for each year

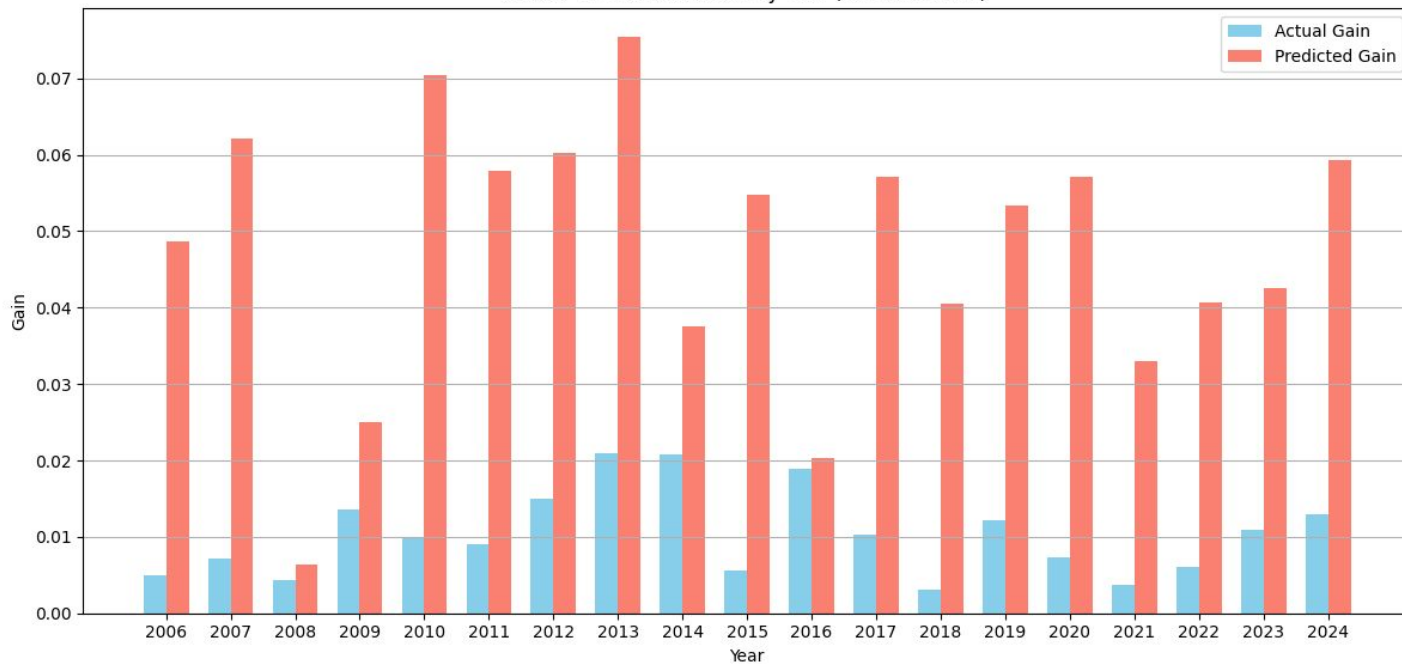
$$\text{pred_gain} = \text{total_entropy} - \text{feature_entropy}$$

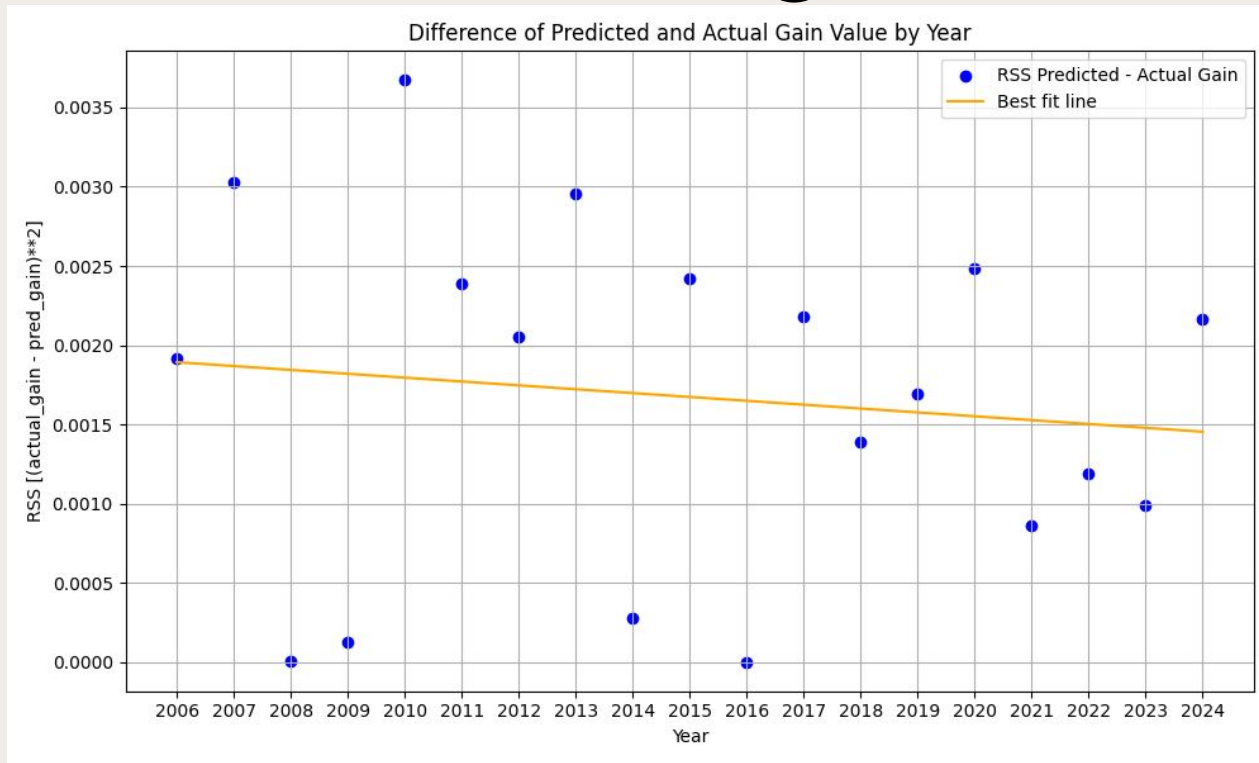
Predicted
value

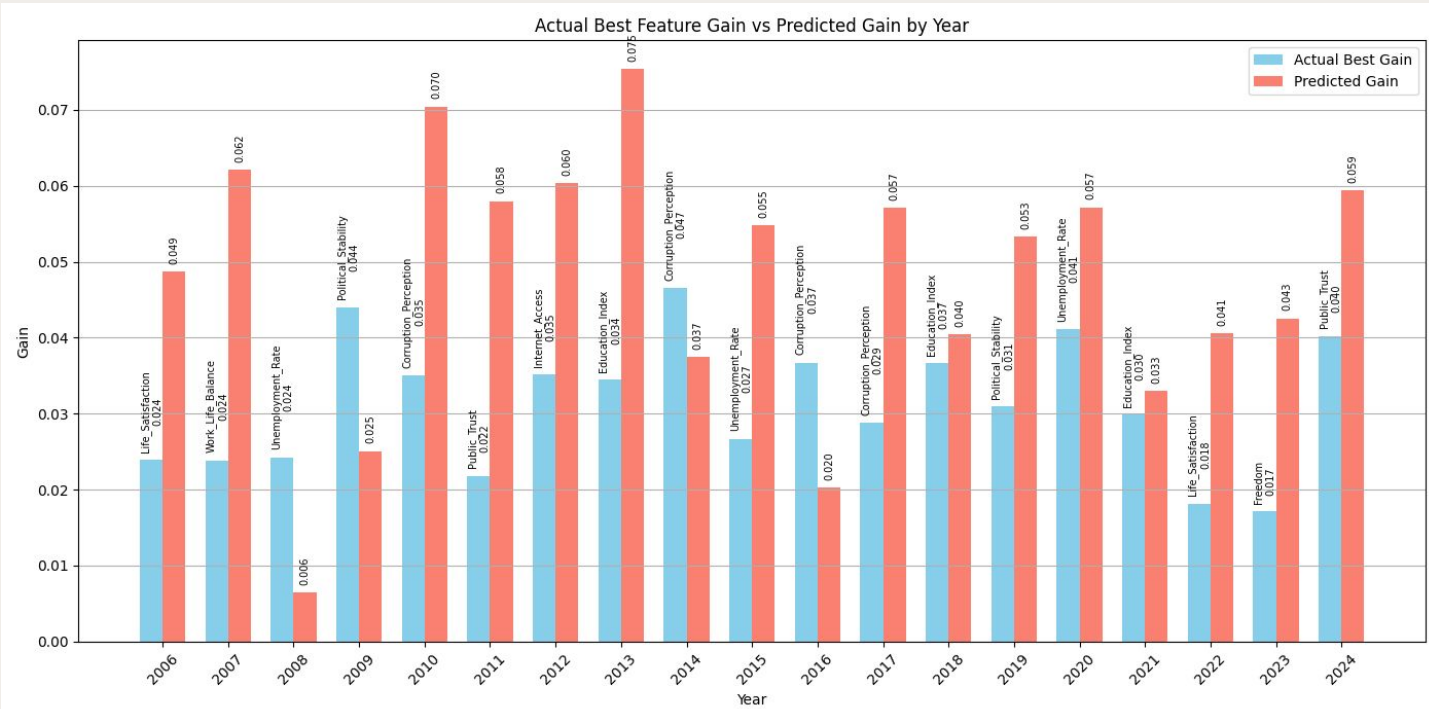
- Calculate for all years

$$\text{difference} = (\text{actual_gain} - \text{pred_gain})^{**2}$$

Actual vs Predicted Gain by Year (Same Feature)









04

Conclusions



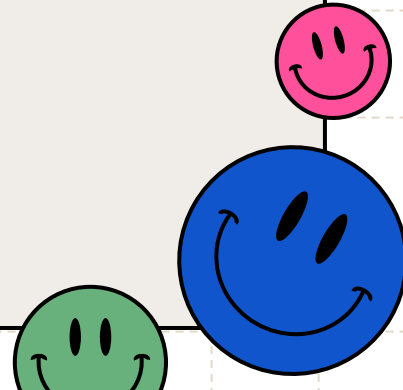
To Summarize...

Main Takeaways

- Public Health Expenditure most directly impact happiness level for **2005 ONLY**
 - Gain values not really high
- Predictions overemphasize the gain value of the feature.
- We **should not** use gain and entropy to **predict** best features.

Future Works

- Implement **bootstrapping or random sampling**
(randomly select predictive data, select smallest RSS iteration)
- Including a **larger dataset**
(more years)



Questions?

