# IS883: Deploying Generative AI

Mohannad Elhamod

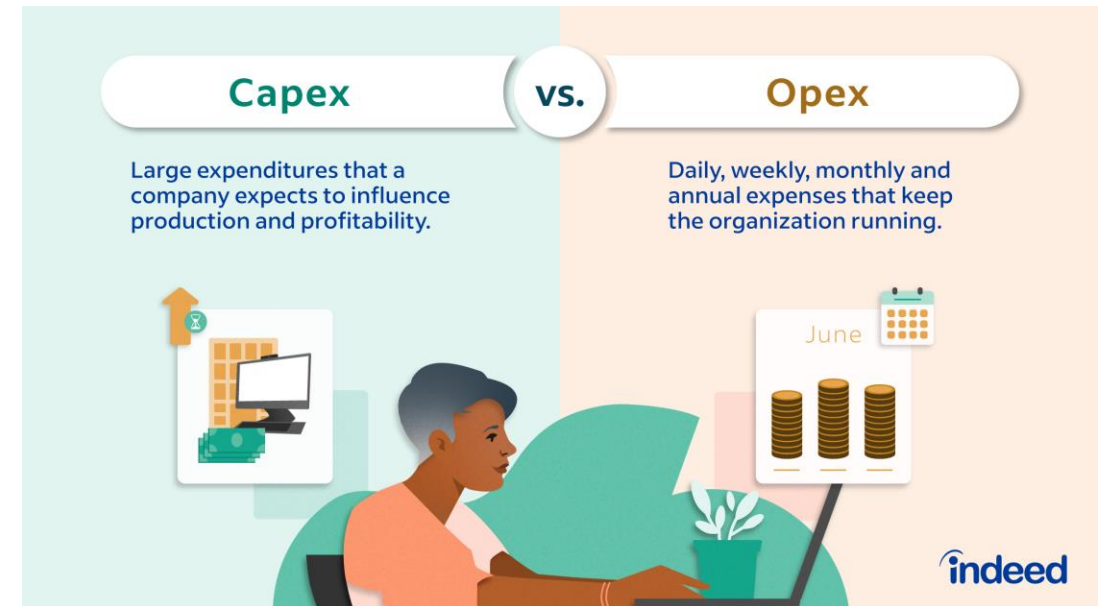**Boston University** Questrom School of Business

# Cost-Benefit Analysis

- When building a Gen AI solution or application, we need to assess its _viability_.

- We generally would look at the _ROI_ (Return on Investment)

$$\text{ROI} = \frac{\text{Net Profit}}{\text{Total Cost}} = \frac{\text{Total Revenue} - \text{Total Cost}}{\text{Total Cost}}$$

# Total Cost

- Total cost can be broken down into two types:
  - **Upfront Costs (CapEx):** Do not repeat.
  - **Ongoing Costs (OpEx):** Recurring expenses to keep the service running.



Capex vs. Opex

**Capex** — Large expenditures that a company expects to influence production and profitability.

**Opex** — Daily, weekly, monthly and annual expenses that keep the organization running.

June

indeed

# Upfront Costs

- Examples


Software and Tools


Data


Infrastructure
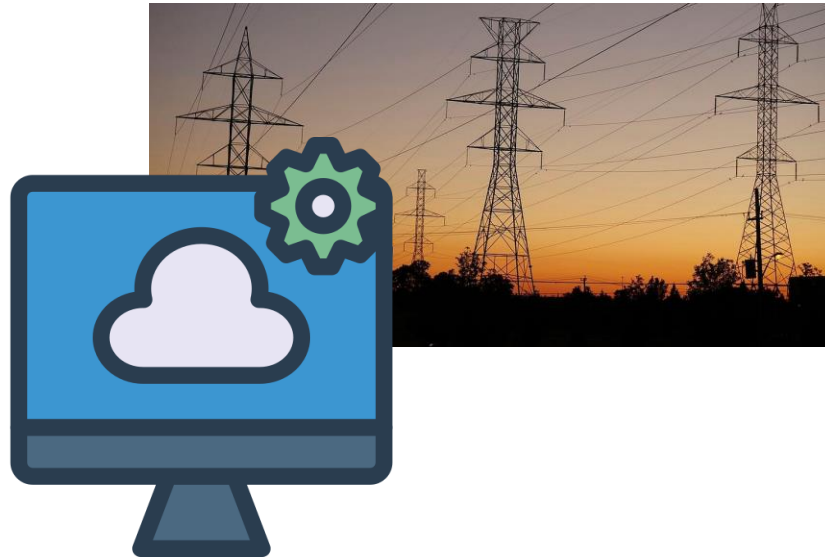

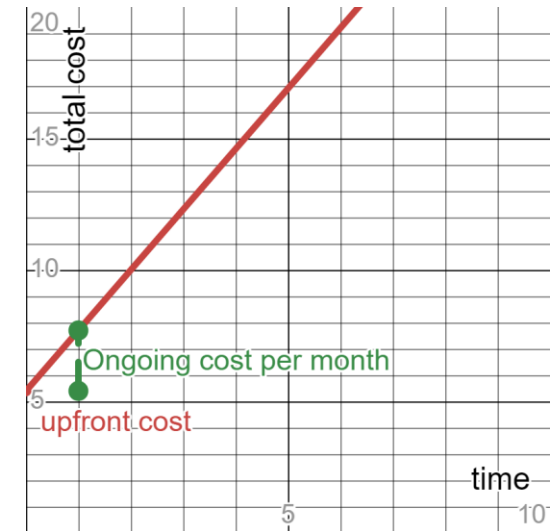Staffing

# Ongoing Costs

- Examples



Operational costs (Energy, subscriptions, etc.)



Maintenance costs

# Why Care For This Distinction?

- Because that's what matters for the cost equation.
  - Your decision should be based on:
    - <u>Cashflow:</u> **How much** money you have and **when**.
    - <u>Product Lifespan.</u>
    - <u>Scalability and Flexibility:</u> Pay-as-you-go provides more control and easier adaptation.
    - <u>Risk:</u> Lower upfront commitment is more risk averse.



$$\text{Total Cost} = \text{Upfront Costs} + (\text{Ongoing Cost per Unit of Time} \times \text{Time Period})$$

# What Affects Cost in Gen AI?

| In-house Approach | Pay-as-you-go Approach |
|---|---|
| • **Infrastructure** (upfront or ongoing).<br>• **Software licensing** (ongoing)<br>• **Developers and Staff** (upfront and ongoing).<br>• **Maintenance, energy, and upgrades** (ongoing).<br>    • Includes model size (computational needs)<br>• **Training cost** (upfront or ongoing):<br>    • Data size | • **Developers and Staff** (upfront and ongoing).<br>• **Cloud subscription fees** (ongoing).<br>• **Usages** (ongoing):<br>    • Number of tokens (i.e., request and response sizes).<br>    • Model type. |

# Total Revenue

- <u>Total revenue</u> depends on the provided service/product, but could be measured as
  - **<u>Cost Savings:</u>** If product or service is to be used internally, how much cost reduction is there due to automation and increased efficiency?
  - **<u>Sales Increase:</u>** due to an increase in customer satisfaction.

$$\text{Total Revenue} = (\text{Sales per Period} \times \text{Price per Unit} + \text{Cost Savings per Period}) \times \text{Time Period}$$

$$\text{Cost Savings} = (\text{Reduction in Personnel} \times \text{Personnel Cost per Person})$$
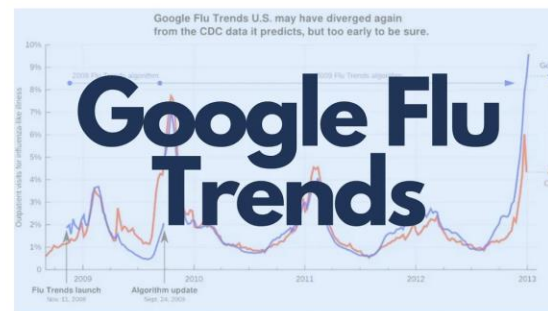
# Evaluations: ROI, CBR, and Break-Even Analysis

- <u>Break-Even Point:</u> Time it takes for Total Revenue to become equal Total Cost.

- <u>Cost-Benefit Ratio</u> is simply:  $\text{CBR} = \text{ROI} + 1$

  - ROI is negative when product/service is not viable.
  - CBR is less than one when product/service is not viable.

- Try it [here](#)!

# Uncertainty and Risk

- May not be straight forward to incorporate.

- Attributed to several factors:
    - Regulations and compliance.
    - Long-term performance.





Why Did Google Flu Trends Fail?

# Risk Evaluation

- We will not discuss this in detail as it requires going into statistics.

- But, in its simplest form, you could account for worst- and best-case scenarios to study the sensitivity of your estimates.

# Resources

- [Gen AI: too much spend, too little benefit?](#)
- [Understanding the Cost of LLMs](#)