# Designing a Machine Learning System -II

# Types of Machine Learning Systems

ML systems can be categorized based on three aspects:
- How their ML models serve their predictions (batch prediction vs. online prediction)
- Where the majority of computation is done (edge computing vs. cloud computing)
- How often their ML models get updated (online learning vs. offline learning)

# Batch prediction vs. online prediction

● Two options for predictions offered by major cloud providers

● Valid when using your own data centers too

● In both cases you pass input data to a cloud-hosted machine-learning model and get inferences for each data instance


Google Cloud

# Batch prediction vs. online prediction

- Batch prediction:
  - Generate predictions periodically
  - Predictions are stored somewhere (e.g. SQL tables, CSV files)
  - Retrieve them as needed
  - Allow more complex models
- Online prediction:
  - Generate predictions as requests arrive
  - Predictions are returned as responses

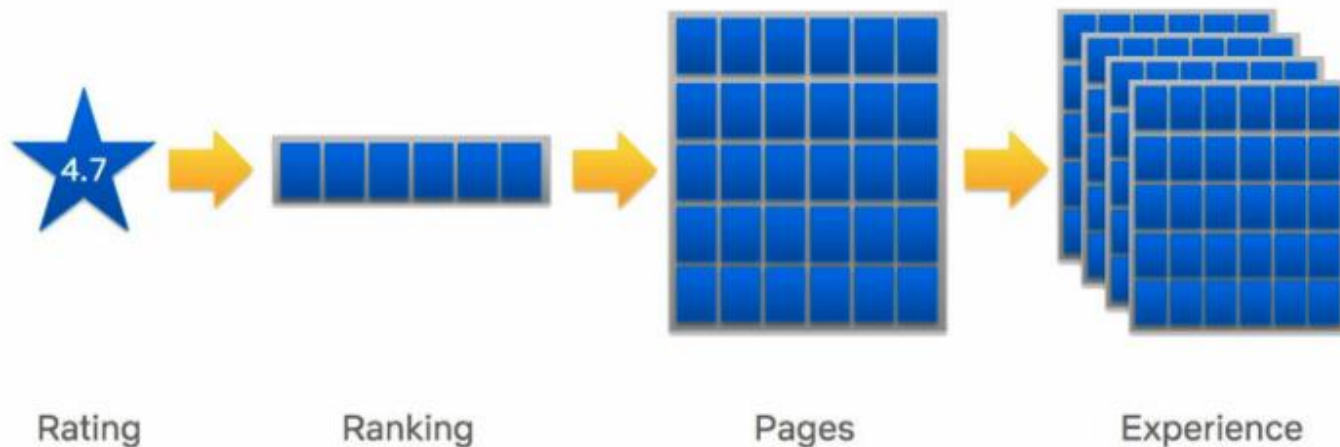|  | **Batch prediction** | **Online prediction** |
|---|---|---|
| **Frequency** | Periodical (e.g. every 4 hours) | As soon as requests come |
| **Useful for** | Processing accumulated data when you don't need immediate results (e.g. recommendation systems) | When predictions are needed as soon as data sample is generated (e.g. fraud detection) |
| **Optimized** | High throughput | Low latency |
| **Input space** | Finite: need to know how many predictions to generate | Can be infinite |
| **Examples** | ● TripAdvisor hotel ranking<br>● Netflix recommendations | ● Google Assistant speech recognition<br>● Twitter feed |

# Netflix

- Netflix is now synonymous to most people as the go-to streaming service for movies and tv shows.
- It has transformed from a mail service posting DVDs in the US (late 90s) to a global streaming service with 182.8 million subscribers.

# Netflix



**Evolution of our Personalization Approach**

Rating → Ranking → Pages → Experience

# Netflix Recommender System



10,000s of possible rows

Variable number of possible videos per row (up to thousands)

1 personalized page

10-40 rows

per device

8

# Netflix Recommender System

Netflix utilizes a two-tiered row-based ranking system, where ranking happens:
- Within each row (strongest recommendations on the left)
- Across rows (strongest recommendations on top)
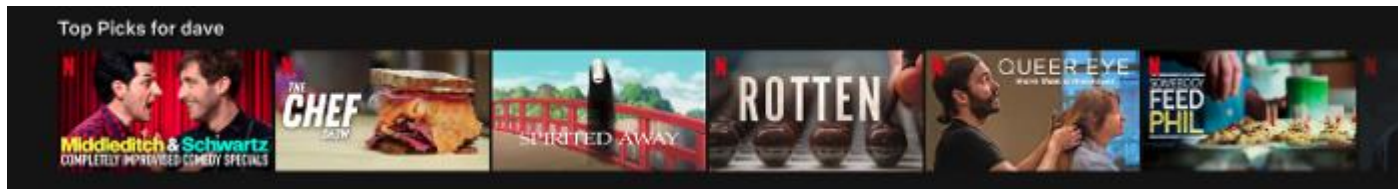
# Netflix – Ranking Algorithms

**Personalised Video Ranking (PVR)** — A general-purpose algorithm, which filters down the catalog by a certain criteria (e.g. Violent TV Programmes, US TV shows, Romance, etc), combined with side features including user features and popularity. Title recommendations use **batch predictions**
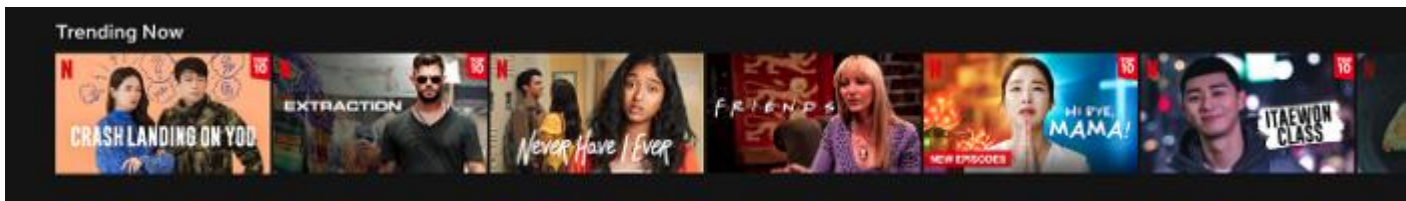
https://towardsdatascience.com/deep-dive-into-netflixs-recommender-system-341806ae3b48

# Netflix – Ranking Algorithms

**Top-N Video Ranker** — Similar to PVR except that it is optimized using metrics that look at the top percentiles of the catalog ranking (e.g. MAP@K, NDCG).

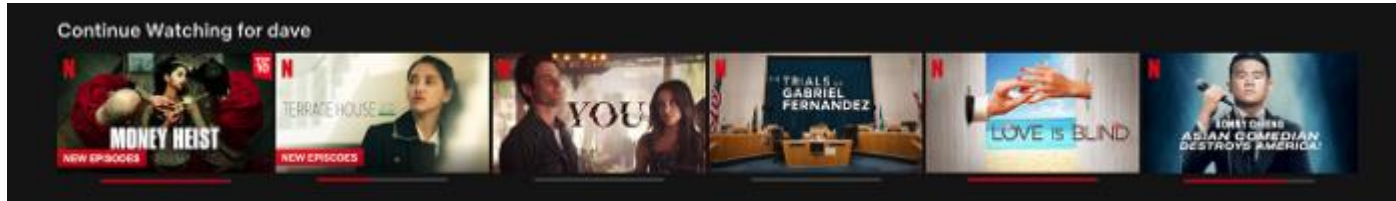https://towardsdatascience.com/deep-dive-into-netflixs-recommender-system-341806ae3b48

# Netflix – Ranking Algorithms

- **Trending Now Ranker** — This algorithm captures temporal trends which Netflix deduces to be strong predictors. These short-term trends can range from a few minutes to a few days.
- These events/trends are typically:
  - Events that have a seasonal trend and repeat themselves (e.g. Valentines day leads to an uptick in Romance videos being consumed)
  - One-off, short term events (e.g. recent disasters leading to short-term interest in documentaries about them)

https://towardsdatascience.com/deep-dive-into-netflixs-recommender-system-341806ae3b48
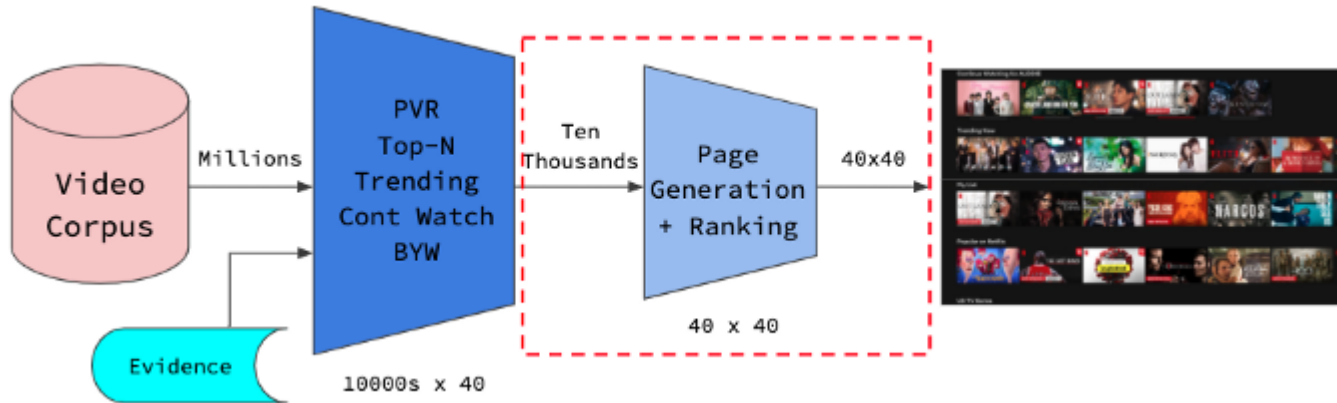
# Netflix – Ranking Algorithms

- **Continue Watching Ranker** — This algorithm looks at items that the member has consumed but has not completed, typically:
- Episodic content (e.g. TV series)
- Non-episodic content that can be consumed in small bites (e.g. movies that are half-completed, series that are episode independent such as Black Mirror)
- The algorithm calculates the probability of the member continue watching and includes other context-aware signals (e.g. time elapsed since viewing, point of abandonment, device watched on).

https://towardsdatascience.com/deep-dive-into-netflixs-recommender-system-341806ae3b48

# Netflix Page Generation Process

- After the algorithms generate candidate rows (already ranked within each row vector), it now needs to decide which of these 10,000s of rows to display?



https://towardsdatascience.com/deep-dive-into-netflixs-recommender-system-341806ae3b48

# Netflix Page Generation Process

- Now there is a set of possible rows to consider for a page, it is time to assemble the homepage from them
- First find candidate groupings that are likely relevant for a member based on the information available for that member.
- This also involves coming up with the evidence (or explanations) to support the presentation of a row, for example the movies that the member has previously watched in a genre.

https://towardsdatascience.com/deep-dive-into-netflixs-recommender-system-341806ae3b48

# Netflix Page Generation Process

- Next, each group is filtered to handle concerns like maturity rating or to remove some previously watched videos.
- After filtering, the videos in each group are ranked according to a row-appropriate ranking algorithm, which produces an ordering of videos such that the most relevant videos for the member in a group are at the front of the row.
- From this set of row candidates a row selection algorithm is run to assemble the full page.

Row selection use **online predictions**

https://towardsdatascience.com/deep-dive-into-netflixs-recommender-system-341806ae3b48

# Edge computing vs. cloud computing

|  | **Cloud computing** | **Edge computing** |
|---|---|---|
| **Computations** | Done on cloud (servers) | Done on edge devices (browsers, phones, tablets, laptops, smart watches, activity watchers, cars, etc.) |
| **Requirements** | **Network connections**: availability and speed for data transfer | **Hardware**: memory, compute power, energy for doing computations |
| **Examples** | Most queries to Alexa, Siri, Google Assistant<br><br>Google Translate for rare language pairs (e.g. English - Yiddish) | Wake words for Alexa, Siri, Google Assistant<br><br>Google Translate for popular language pairs (e.g. English - Spanish)<br><br>Predictive text<br><br>Unlocking with fingerprints, faces |

# Benefits of edge computing

- Can work without (Internet) connections or with unreliable connections
  - Many companies have strict no-Internet policy
  - **Caveat**: devices are capable of doing computations but apps need external information
    - e.g. ETA needs external real-time traffic information to work well
- Don't have to worry about network latency
  - Network latency might be a bigger problem than inference latency
  - Many use cases are impossible with network latency
    - e.g. predictive texting
- Fewer concerns about privacy
  - Don't have to send user data over networks (which can be intercepted)
  - Cloud database breaches can affect many people
  - Easier to comply with regulations (e.g. GDPR)
  - **Caveat**: edge computing might make it easier to steal user data by just taking the device
- Cheaper
  - The more computations we can push to the edge, the less we have to pay for servers

# Offline learning vs. Online learning

| | Offline learning | Online learning |
|---|---|---|
| **Iteration cycle** | Periodical (months) | Continual (minutes) ≠ continuous |
| **Batch size** | batch (thousands -> millions of samples)<br>GPT-3 125M params: batch size 0.5M<br>GPT-3 175B params: batch size 3.2M | microbatch (hundreds of samples) |
| **Data usage** | Each sample seen multiple times (epochs) | Each sample seen at most once |
| **Evaluation** | Mostly offline evaluation | Offline evaluation as sanity check<br>Mostly relying on online evaluation (A/B testing) |
| **Examples** | Most applications | TikTok recommendation system<br><br>Twitter hashtag trending |

# Online learning vs. offline learning

- Both can be done together to create more stable systems
- **If the infrastructure is set up right**, there's no fundamental difference between online learning and offline learning, just a hyperparameter to tune.

# ML in production: expectation

1. Collect data
2. Train model
3. Deploy model



Collect data

Train model

Deploy model

**Waterfall model**

# ML in production: reality

1.  Choose a metric to optimize
2.  Collect data
3.  Train model
4.  Realize many labels are wrong -> relabel data
5.  Train model
6.  Model performs poorly on one class -> collect more data for that class
7.  Train model
8.  Model performs poorly on most recent data -> collect more recent data
9.  Train model
10. Deploy model
11. There are complaints that model biases against one group! -> revert to older version
12. Get more data, train more, do more testing
13. Deploy model
14. Model performs well but revenue decreasing
15. Choose a different metric
16. Start over

# ML in production: reality
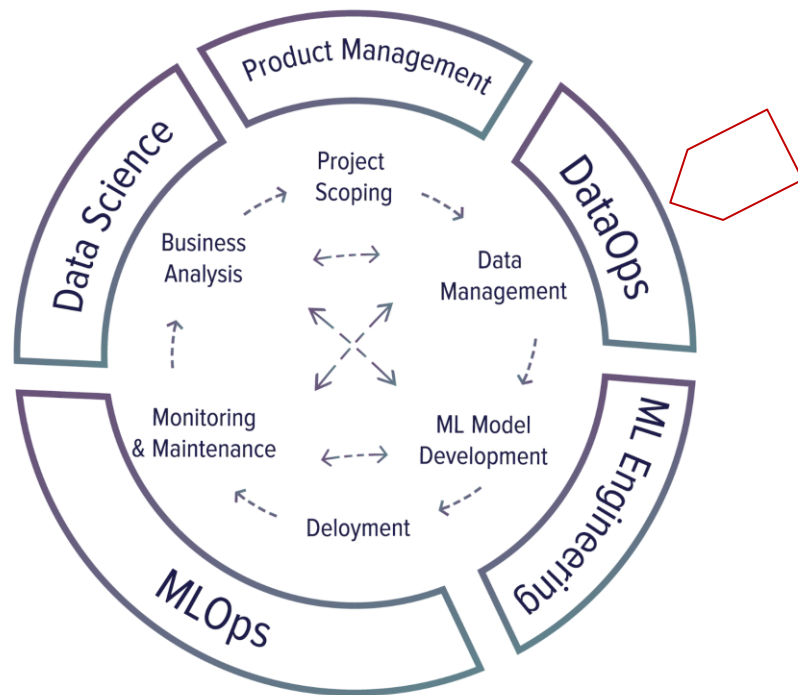
1. Choose a metric to optimize
2. Collect data
3. Train model
4. Realize many labels are wrong -> relabel data
5. Train model
6. Model performs poorly on one class -> collect more data for
7. Train model
8. Model performs poorly on most recent data -> collect more
9. Train model
10. Deploy model
11. There are complaints that model biases against one group!
12. Get more data, train more, do more testing
13. Deploy model
14. Model performs well but revenue decreasing
15. Choose a different metric
16. Start over

**Iterative development**

# Project scoping

- Goals & objectives
- Constraints
- Evaluation

# Data management

- Data sources
- Data format
- Processing
- Storage
- Data consumer
- Data controller

# Model development

- Dataset creation
- Feature engineering
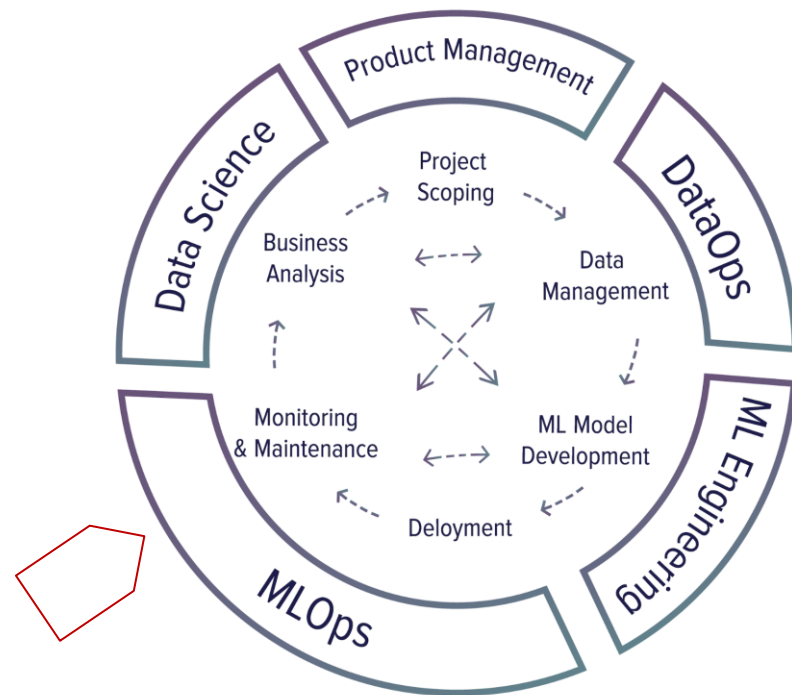- Model training
- Offline model evaluation

# Deployment

- Deploying and serving
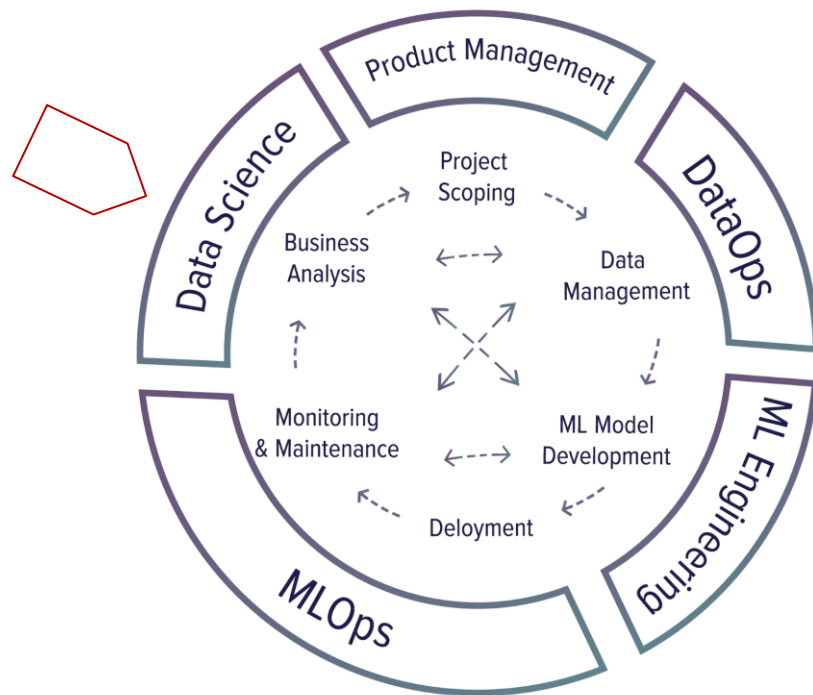- Release strategies
- Online model evaluation

# Monitoring & maintenance

- Model performance & data monitoring
- Model retraining
- Model updates

# Business analysis

- User experience
- Tying model performance to business performance

# Project scoping
## Goals

The ultimate goal of an ML project or any project within a business is, to increase profits directly or indirectly.

- Directly: increasing sales (ads, conversion rates), cutting costs
- Indirectly: increasing customer satisfaction, increasing time spent on a website

**Non-profits are exceptions**: lots of exciting applications of AI for social good
- environment (climate change, deforestation, flood risk, etc.)
- public health
- education (intelligent tutoring system, personalized learning)

# ML to business performance can be confusing

An ML model that gives customers more personalized solutions can either:

- make them happier which makes them spend more money
- solve their problems faster which makes them spend less money

# Example

Possible goals when building a ranking system for newsfeed?

1. Minimize the spread of misinformation
2. Maximize revenue from sponsored content
3. Maximize engagement

# Objectives

| Goals | Objectives |
|---|---|
| General purpose of a project | Specific steps on how to realize that purpose |

# Example: ranking system for newsfeed

| Goals | Objectives |
|---|---|
| General purpose of a project | Specific steps on how to realize that purpose |
| Maximize users' engagement | 1. Filter out spam<br>2. Filter out NSFW content<br>3. Rank posts by engagement: how likely users will click on it |

# Example: ranking system for <u>wholesome</u> newsfeed

| Goals | Objectives |
|---|---|
| General purpose of a project | Specific steps on how to realize that purpose |
| **Maximize users' engagement** while **minimizing the spread of extreme views** and **misinformation** | 1. Filter out spam<br>2. Filter out NSFW content<br>3. Filter out misinformation<br>4. Rank posts by quality<br>5. Rank posts by engagement: how likely users will click on it |

# Example: ranking system for wholesome newsfeed

| Goals | Objectives |
|---|---|
| General purpose of a project | Specific steps on how to realize that purpose |
| **Maximize users' engagement** while **minimizing the spread of extreme views** and **misinformation** | 1. Filter out spam<br>2. Filter out NSFW content<br>3. Filter out misinformation<br>4. Rank posts by quality<br>5. Rank posts by how likely users will click on it |

How to combine these two objectives?

# Multiple objective optimization (MOO)

- Rank posts by quality
  - Predict posts' quality
  - Minimize **quality_loss:** difference between predicted quality and true quality

- Rank posts by how likely users will click on it
  - Predict posts' engagement
  - Minimize **engagement_loss**: difference between predicted clicks and true clicks

# MOO: one model with combined loss

- Rank posts by quality
  - Predict posts' quality
  - Minimize **quality_loss:** difference between predicted quality and true quality

- Rank posts by how likely users will click on it
  - Predict posts' engagement
  - Minimize **engagement_loss**: difference between predicted clicks and true clicks

$$loss = \alpha \text{ quality\_loss} + \beta \text{ engagement\_loss}$$

Train one model to minimize this combined loss
Tune $\alpha$ and $\beta$ to meet your need

# MOO: one model with combined loss

- Rank posts by quality
  - Predict posts' quality
  - Minimize **quality_loss:** difference between predicted quality and true quality

- Rank posts by how likely users will click on it
  - Predict posts' engagement
  - Minimize **engagement_loss**: difference between predicted clicks and true clicks

$$\text{loss} = \alpha \text{ quality\_loss} + \beta \text{ engagement\_loss}$$

Every time you want to tweak $\alpha$ and $\beta$, you have to retrain your model!

# MOO: combine different models

- Rank posts by quality
  - Predict posts' quality
  - Minimize **quality_loss:** difference between predicted quality and true quality

- Rank posts by how likely users will click on it
  - Predict posts' engagement
  - Minimize **engagement_loss**: difference between predicted clicks and true clicks

$M_q$: optimizes **quality_loss**
$M_e$: optimizes **engagement_loss**

Rank posts by $\alpha$ **$M_q$(post)** + $\beta$ **$M_e$(post)**

Now you can tweak $\alpha$ and $\beta$ without retraining models

# Decouple different objectives

- Easier for training:
  - Optimizing for one objective is easier than optimizing for multiple objectives
- Easier to tweak your system:
  - E.g. $\alpha$ % model optimized for quality + $\beta$ % model optimized for engagement
- Easier for maintenance:
  - Different objectives might need different maintenance schedules
    - Spamming techniques evolve much faster than the way post quality is perceived
    - Spam filtering systems need updates more frequently than quality ranking systems

# Netflix Prize

- In 2000, Netflix introduced personalized movie recommendations and in 2006, launched **Netflix Prize,** a machine learning and data mining competition with a $1 million dollar prize money.
- Back then, Netflix used *Cinematch*, its proprietary recommender system which had a *root mean squared error* (RMSE) of 0.9525 and challenged people to beat this benchmark by 10%.
- The team who could achieve the target or got close to this target after a year would be awarded the prize money.

# Netflix Prize - 2007

- The winner in 2007 used a linear combination of *Matrix Factorization* and *Restricted Boltzmann Machines*, achieving a RMSE of 0.88. Netflix then put those algorithms into production after some adaptations.

# Netflix Prize - 2009

● Despite some teams achieving a RMSE of 0.86 in 2009, the company did not put those algorithms into production due to the engineering effort required to gain the marginal increase in accuracy.
   An important point in real-life recommender systems — that there is always a positive relationship between model improvements and engineering efforts.

# When to use ML

Machine learning is an approach to <u>learn</u> <u>complex patterns</u> from <u>existing data</u> and use these patterns to make <u>predictions</u> on <u>unseen data</u>.

# When not to use ML

- It is unethical.
- Simpler solutions do the trick.
- It is impossible to get the right data.
- One single prediction error can cause devastating consequences.
- Every single decision the system makes must be explainable.
- It's not cost-effective.

# If ML is nott the solution, it can be part of it

● Break your problem into smaller problems. ML might be able to solve some of them.
  ○ If a chatbot can't answer all customers' queries, build an ML model to predict whether a query matches one of the frequently asked questions.
  ○ Human-in-the-loop

# Four phases of ML adoption

## Phase 1: Before ML

"If you think that machine learning will give you a 100% boost, then a heuristic will get you 50% of the way there."

Martin Zinkevich, Google

https://newsfeed.org/what-mark-zuckerbergs-news-feed-looked-like-in-2006/

# Four phases of ML adoption

## Phase 2: Simplest ML models

Start with a simple model that allows visibility into its working to:

- validate hypothesis
- validate pipeline

# Four phases of ML adoption

## Phase 3: Optimizing simple models

- Different objective functions
- Feature engineering
- More data
- Ensembling

# Four phases of ML adoption

## Phase 4: Complex ML models