

## Quiz 2 - 563

### Word Embeddings

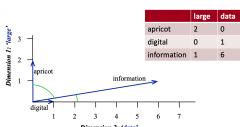
- Motivation:** You can understand a word by the context/company it keeps.

### Introduction

- Standard approach: put words in vector space and the distance between words is the similarity between them.
- word2vec is unsupervised/ semi-supervised learning because:
  - closely related to dimensionality reduction + extracting meaningful representation from raw data
  - do not need any labeled data
  - running text is used as supervision signal

### Word Representation

- One-hot representation:**
  - Simplest way to represent a word
  - OHE vector is a vector of all 0s except for a 1 at the index of the word in the vocabulary
  - rows = words in sentence, columns = words in vocabulary
- Disadvantages:**
  - High dimensionality
  - No notion of similarity between words (dot product = 0)
  - No notion of context
- Term-term co-occurrence matrix:**
  - A matrix where each row and column corresponds to a word in the vocabulary
  - The value in the  $i$ -th row and  $j$ -th column is the number of times word  $i$  and word  $j$  appear together in a context window
  - Context window:** a fixed-size window that slides over the text (e.g. window size = 2 means 2 words to the left and 2 words to the right)
- Disadvantages:**
  - High dimensionality
  - Sparse
  - Does not capture polysemy (multiple meanings of a word)



### Dense Word Representations

- Term-term co-occurrence matrix is sparse and high-dimensional
- Better to learn short and dense vectors for words
  - Easier to store and train
  - Better generalization
- Approaches:
  - Latent Semantic Analysis (LSA):** Use SVD to reduce the dimensionality of the term-term co-occurrence matrix
    - Works better for small datasets compared to word2vec
  - Word2Vec:** Use neural networks to learn word embeddings

### Word2Vec

- Create short and dense word embeddings using neural networks
- Ideas: Predict the context of a word given the word itself
- skip-gram: Predict the context words given the target word
- Continuous Bag of Words (CBOW): Predict the target word given the context words
- Two moderately efficient training algorithms:
  - Hierarchical softmax: Use a binary tree to represent all words in the vocabulary
  - Negative sampling: Treat the problem as a binary classification problem

### word2vec: Skip-gram

- Predict the context words given the target word
- NN to obtain short and dense word vectors
- Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

• NN to obtain short and dense word vectors

• Architecture:

word2vec: Skip-gram

• Predict the context words given the target word

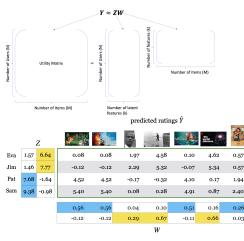
• NN to obtain short and dense word vectors

• Architecture:

## Collaborative Filtering

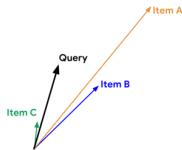
- Unsupervised learning
- **Intuition:**
  - People who agreed in the past are likely to agree again in future
  - Leverage social information for recommendations
- **PCA ?:**
  - To learn latent features of users and items
  - Run on utility matrix
  - **Problem:** missing values
    - PCA loss function  $J(Z, W) = \sum_{i,j} ||W^T Z_{ij} - Y_{ij}||^2$
    - Cannot use SVD directly because many missing values AND missing values make SVD undefined
  - **Solutions:**
    - Impute the values to do PCA
      - BUT, will introduce bias (distort the data)
      - Bias will be compounded by the imputed values
    - Summing over only available values
    - Prone to overfitting
  - **Collaborative Filtering Loss Function:**
    - Only consider the available values
    - Add L2-reg to the loss function for  $W$  and  $Z$
    - $J(Z, W) = \sum_{i,j} ||W^T Z_{ij} - Y_{ij}||^2 + \frac{\lambda}{2} ||W||^2 + \frac{\lambda}{2} ||Z||^2$
    - This accounts for the missing values and the regularization terms prevent overfitting (representations are not too complex)
    - This improved the RMSE score by 7% in the Netflix competition
    - Optimize using SGD (stochastic gradient descent) and WALS (weighted alternating least squares)
- **Other Notes:**
  - Result can be outside the range of the ratings
  - Will have problems with cold start (new users or items)
  - Use `surprise` library for collaborative filtering

## Z and W in Collaborative Filtering



- $Z$  is no longer the points in the new hyperplane and  $W$  is no longer the weights
- $Z$ : row = user | Maps users to latent feature of documents
  - e.g. left for animation, right for documentaries
- $W$ : col = item | Maps items to latent feature of users
  - e.g. top for animation lovers, bottom for documentary lovers

## Distance Metrics



- **Cosine:**
  - $d(x, y) = \frac{x \cdot y}{\|x\| \|y\|}$
  - Collinear = 1, orthogonal = 0 (want a value close to 1)
  - It is the angle between the two vectors
  - Rank (high to low): C, A, B
- **Euclidean:**
  - $d(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$
  - It is the straight line distance between the two points (want smaller distance)
  - Rank (low to high): B, C, A
- **Dot Product:**
  - $d(x, y) = x \cdot y$
  - It is the projection of one vector onto the other
  - If vectors are normalized, it is the same as cosine similarity (want larger value)
  - Rank (high to low): A, B, C

## Content Based Filtering

- Supervised learning
- Does not make use of social network / information
- Solves the cold start problem (can recommend items to new users/items)
- Assumes that we have **features of items and/or users** to predict ratings
- Create a user profile for each user
  - Treat rating prediction as a regression problem
  - Have a regression model for each user

## Steps in Python (With a movie recommendation example)

0. Load `ratings_df` (contains `user_id`, `movie_id`, and `rating`)
  - Also make the user and movie mappers.
1. Load movie features. This is a matrix of shape `(n_movies, n_features)`
  - Index of movie features is movie id/name
  - Features can be genre, director, actors, etc.

|                     | Action | Romance | Drama | Comedy | Children | Documentary |
|---------------------|--------|---------|-------|--------|----------|-------------|
| A Beautiful Mind    | 0      | 1       | 1     | 0      | 0        | 0           |
| Bambi               | 0      | 0       | 1     | 0      | 1        | 0           |
| Cast Away           | 0      | 1       | 1     | 0      | 0        | 0           |
| Die Hard            | 0      | 0       | 0     | 0      | 0        | 1           |
| Inception           | 1      | 0       | 1     | 0      | 0        | 0           |
| Jerry Maguire       | 0      | 1       | 1     | 1      | 0        | 0           |
| Lion King           | 0      | 0       | 1     | 0      | 1        | 0           |
| Mission: Impossible | 0      | 0       | 0     | 0      | 0        | 1           |
| Man on Wire         | 0      | 0       | 0     | 0      | 0        | 1           |
| Roman Holidays      | 0      | 1       | 1     | 1      | 0        | 0           |
| The Social Network  | 0      | 0       | 0     | 0      | 0        | 1           |
| Titanic             | 0      | 1       | 1     | 0      | 0        | 0           |

2. Build a user profile. For each user, we will get the ratings and the corresponding movie features.

- Results in a dictionary (key: user, value: numpy array of size `(n_ratings, n_genres)`)

- `n_ratings`: number of movies rated by the user
- 3. **Supervised learning:** Train a regression model for each user.
  - We will use `Ridge` regression model for this example.
  - $x$  is movie features (shape = `(n_movies, n_features)`)
  - $y$  is the ratings (shape = `(n_movies, 1)`)
- Results in "known" weights for each feature for each user.  $W$  of size `(n_features, 1)` for each user.
 

|             | Coefficients for Eva |  |  |  |  |  |
|-------------|----------------------|--|--|--|--|--|
| Action      | 0.333333             |  |  |  |  |  |
| Romance     | -0.000000            |  |  |  |  |  |
| Drama       | -0.666667            |  |  |  |  |  |
| Comedy      | -0.666667            |  |  |  |  |  |
| Children    | 0.000000             |  |  |  |  |  |
| Documentary | 0.666667             |  |  |  |  |  |

## Collaborative vs Content Based Filtering

| Collaborative Filtering                       | Content Based Filtering  |
|---|--|
| $\hat{y}_{ij} = w_j^T z_i$                    | $\hat{y}_{ij} = w_i^T x_{ij}$  |
| $w_j^T$ : "hidden" embedding for feature $j$  | $w_i$ : feature vector for user $i$  |
| $z_i$ : "hidden" embedding for user $i$       | $x_{ij}$ : feature $j$ for user $i$  |
| (+) Makes use of social network / information | (-) Does not make use of social network / information                                  |
| (-) Cold start problem                        | (+) Solves the cold start problem  |
| (+) No feature engineering required           | (-) Requires feature engineering   |
| (-) Hard to interpret                         | (+) Easy to interpret  |
| (-) Cannot capture unique user preferences    | (+) Can capture unique user preferences (since each model is unique)                   |
| (+) More diverse recommendations              | (-) Less diverse recommendations (hardly recommend an item outside the user's profile) |

## Beyond Error Rate in Recommender Systems

- Best RMSE ≠ Best Recommender
- Need to consider simplicity, interpretability, code maintainability, etc.
  - The Netflix Prize: The winning solution was never implemented

- Other things to consider:
  - Diversity: If someone buys a tennis racket, they might not want to buy another tennis racket
  - Freshness: New items (new items need to be recommended for it to be successful)
  - Trust: Explain your recommendation to the user
  - Persistence: If the same recommendation is made over and over again, it might not be a good recommendation
  - Social influence: If a friend buys something, you might want to buy it too
- Also need to consider ethical implications
  - Filter bubbles
  - Recommending harmful content

$$\bullet d(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

- It is the straight line distance between the two points (want smaller distance)
- Rank (low to high): B, C, A

## Dot Product

- $d(x, y) = x \cdot y$
- It is the projection of one vector onto the other
- If vectors are normalized, it is the same as cosine similarity (want larger value)
- Rank (high to low): A, B, C

## Content Based Filtering

- Supervised learning
- Does not make use of social network / information
- Solves the cold start problem (can recommend items to new users/items)
- Assumes that we have **features of items and/or users** to predict ratings
- Create a user profile for each user
  - Treat rating prediction as a regression problem
  - Have a regression model for each user

## Steps in Python (With a movie recommendation example)

0. Load `ratings_df` (contains `user_id`, `movie_id`, and `rating`)

- Also make the user and movie mappers.

1. Load movie features. This is a matrix of shape `(n_movies, n_features)`

- Index of movie features is movie id/name

- Features can be genre, director, actors, etc.

|                     | Action | Romance | Drama | Comedy | Children | Documentary |
|---------------------|--------|---------|-------|--------|----------|-------------|
| A Beautiful Mind    | 0      | 1       | 1     | 0      | 0        | 0           |
| Bambi               | 0      | 0       | 1     | 0      | 1        | 0           |
| Cast Away           | 0      | 1       | 1     | 0      | 0        | 0           |
| Die Hard            | 0      | 0       | 0     | 0      | 0        | 1           |
| Inception           | 1      | 0       | 1     | 0      | 0        | 0           |
| Jerry Maguire       | 0      | 1       | 1     | 1      | 0        | 0           |
| Lion King           | 0      | 0       | 1     | 0      | 1        | 0           |
| Mission: Impossible | 0      | 0       | 0     | 0      | 0        | 1           |
| Man on Wire         | 0      | 0       | 0     | 0      | 0        | 1           |
| Roman Holidays      | 0      | 1       | 1     | 1      | 0        | 0           |
| The Social Network  | 0      | 0       | 0     | 0      | 0        | 1           |
| Titanic             | 0      | 1       | 1     | 0      | 0        | 0           |

2. Build a user profile. For each user, we will get the ratings and the corresponding movie features.

- Results in a dictionary (key: user, value: numpy array of size `(n_ratings, n_genres)`)