

Recommendation systems

Christos Dimitrakakis

October 23, 2019

Recommendation systems

- Least squares representation

- Preferences as a latent variable

- The recommendation problem

More fun with latent variable models

Social networks

Sequential structures



The recommendation problem

At time t

1. A customer x_t appears.
2. We present a choice a_t .
3. The customer chooses y_t .
4. We obtain a reward $r_t = \rho(a_t, y_t) \in \mathbb{R}$.







The two problems in recommendation systems

- ▶ The modelling (or prediction) problem.
- ▶ The recommendation problem.

How to predict user preferences?

Example: Item-based CF



	2			4	5	2.94*
	5		4			1
			5		2	2.48*
		1		5		4
			4			2
	4	5		1		1.12*

sim(i,j) -1 -1 0.86 1 NA

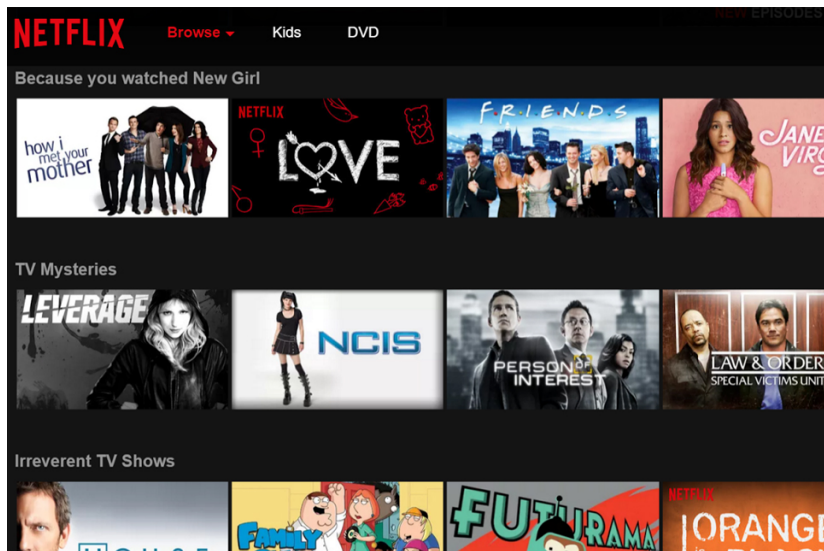


Figure: What to recommend?

Predictions based on similarity

Content-based filtering.

- ▶ Users typically like similar items.
- ▶ That means we can use one user's ratings and **item information** to predict their ratings for other items.

Collaborative filtering

- ▶ **Similar users have similar tastes.**
- ▶ That means we can use similar user's **ratings** to predict the ratings for other users.

k -NN for similarity

Exercise 1

- ▶ Define a distance $d : \mathcal{X}^M \times \mathcal{X}^M \rightarrow \mathbb{R}_+$ between user ratings.
- ▶ Apply a k -NN-like algorithm to prediction of user ratings from the dataset.

Similarity between users

$$\sum_{j \neq i} w_{i,j} = 1, \quad w_{i,j}^m \triangleq w_{i,j} \mathbb{I}\{x_{j,m}\} / \sum_k w_{i,k} \mathbb{I}\{x_{k,m}\}.$$

Example 1 (k -nearest neighbours)

$w_{i,j} = 1/k$ for the k nearest neighbours with respect to d .

Example 2 (Weighted distance)

$$w_{i,j} = \frac{\exp[-d(i,j)]}{\sum_{k \neq i} \exp[-d(i,k)]}$$

Inferred ratings

$$\hat{x}_{u,m} = \sum_{j \neq u} w_{u,j}^m x_{j,m}.$$

A naive distance metric

$$d(i, j) \triangleq \|\mathbf{x}_i - \mathbf{x}_j\|_1.$$

Ignoring movies which are not shared.

$$d(i, j) \triangleq \sum_m \mathbb{I}\{\mathbf{x}_{i,m} \wedge \mathbf{x}_{j,m}\} |\mathbf{x}_{i,m} - \mathbf{x}_{j,m}|$$

Using side-information

Social network data

Inferring a latent representation

$$d(i, j) \triangleq f(\mathbf{x}_i, \mathbf{x}_j, \theta)$$

Latent representation

The predictive model

- ▶ x_{um} rating of user u for movie m .
- ▶ $r_{um} = \mathbb{I}\{x_{um} > 0\}$ indicates which movies are rated.
- ▶ $\mathbf{z}_m \in \mathbb{R}^n$: an n -dimensional representation of a movie.
- ▶ $\mathbf{c}_u \in \mathbb{R}^n$: an n -dimensional representation of a user.

Given \mathbf{C}, \mathbf{Z} , our predicted movie rating can be written as

$$\hat{x}_{u,m} \triangleq \mathbf{c}_u^\top \mathbf{z}_m, \quad \hat{\mathbf{X}} \triangleq \mathbf{C}^\top \mathbf{Z}.$$

$$f(\mathbf{C}, \mathbf{Z}) = \|(\mathbf{R} \circ \hat{\mathbf{X}} - \mathbf{R} \circ \mathbf{X})^\top (\mathbf{R} \circ \hat{\mathbf{X}} - \mathbf{R} \circ \mathbf{X})\|_1$$

A simple preference model

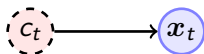


Figure: Basic preference model

Example 3 (Discrete preference model)

- ▶ User type $c \in \mathcal{C}$.
- ▶ User ratings x with $x_m \in \mathcal{X} = \{0, 1\}$ rating for movie m .
- ▶ Preference distribution

$$P_{\theta}(x|c) = \prod_{m=1}^M \theta_{m,c}^{x_m} (1 - \theta_{m,c})^{(1-x_m)}.$$

- ▶ $P_{\theta}(c) = \theta_c, \sum_c \theta_c = 1.$

A simple preference model

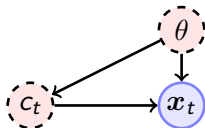


Figure: Basic preference model

Example 3 (Discrete preference model)

- ▶ User type $c \in \mathcal{C}$.
- ▶ User ratings x with $x_m \in \mathcal{X} = \{0, 1\}$ rating for movie m .
- ▶ Preference distribution

$$P_{\theta}(x|c) = \prod_{m=1}^M \theta_{m,c}^{x_m} (1 - \theta_{m,c})^{(1-x_m)}.$$

- ▶ $P_{\theta}(c) = \theta_c, \sum_c \theta_c = 1.$

A more complex preference model

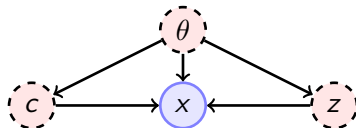


Figure: Preference model

Preference model

- ▶ User type $c \in \mathcal{C}$.
- ▶ Movie type $z \in \mathcal{Z}$.
- ▶ Preference distribution

$$P_{\theta}(x|c, z) = \mathcal{N}(c^{\top} z, \sigma_{\theta})$$

- ▶ Feature prior

$$P_{\theta}(c) = \mathcal{N}(0, \lambda_{\theta})$$

What to recommend

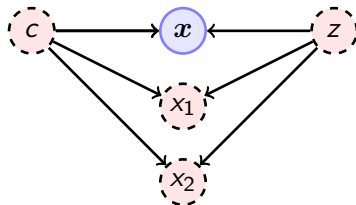


Figure: Preference model

The recommendation problem for a given θ

What to recommend

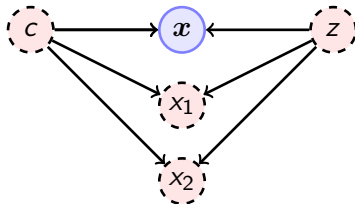


Figure: Preference model

The recommendation problem for a given θ

$$\max_{\pi} \mathbb{E}_{\theta}^{\pi}(U \mid \mathbf{x}) = \max_a \sum_{c,z} U(a, y) \mathbb{P}(y \mid \mathbf{a}, c, z) P_{\theta}(c, z \mid \mathbf{x}) \quad (1.1)$$

$$= \max_a \sum_{c,z} U(a, y) \sum_{x_a} \mathbb{P}(y \mid \mathbf{a}, x_a) P_{\theta}(x_a \mid c, z) P_{\theta}(c, z \mid \mathbf{x}) \quad (1.2)$$

Two ways to model the effect of actions

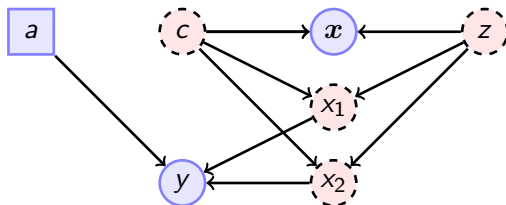


Figure: Preference model

$$\mathbb{E}_{\theta}(U \mid a, x) = \sum_{c,z} U(a, y) \sum_{x_a} \mathbb{P}(y \mid a, x_a) P_{\theta}(x_a \mid c, z) P_{\theta}(c, z \mid x) \quad (1.3)$$

Two ways to model the effect of actions

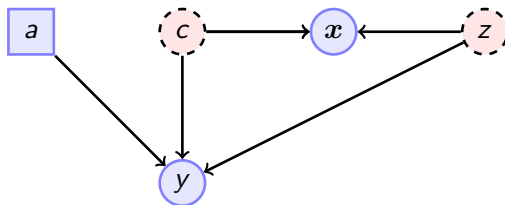


Figure: Preference model

$$\mathbb{E}_{\theta}(U \mid a, x) = \sum_{c, z} U(a, y) \mathbb{P}(y \mid a, c, z) P_{\theta}(c, z \mid x) \quad (1.3)$$

Recommendation systems

More fun with latent variable models

Social networks

Sequential structures

Clusters as latent variables

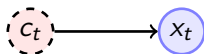


Figure: Graphical model for independent data from a cluster distribution.

The clustering distribution

- ▶ Cluster c_t
- ▶ Observation x_t
- ▶ Parameter θ .

$$x_t \mid c_t = c, \theta \sim P_\theta(x|c), \quad c_t \mid \theta \sim P_\theta(c), \quad \theta \sim \xi(\theta)$$

$$P_\theta(c_t \mid x_t) =$$

Clusters as latent variables

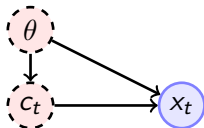


Figure: Graphical model for independent data from a cluster distribution.

The clustering distribution

- ▶ Cluster c_t
- ▶ Observation x_t
- ▶ Parameter θ .

$$x_t \mid c_t = c, \theta \sim P_\theta(x|c), \quad c_t \mid \theta \sim P_\theta(c), \quad \theta \sim \xi(\theta)$$

$$P_\theta(c_t \mid x_t) = \frac{P_\theta(x_t \mid c_t)P_\theta(c_t)}{\sum_{c'} P_\theta(x_t \mid c_t = c')P_\theta(c_t = c')}$$

Bayesian formulation of the clustering problem

- ▶ Prior ξ on parameter space Θ .
- ▶ Data $x^T = x_1, \dots, x_T$. Cluster assignments c^T unknown.
- ▶ Posterior $\xi(\cdot | x^T)$.

Posterior distribution

$$\xi(\theta | x^T) = \frac{P_\theta(x^T)\xi(\theta)}{\sum_{\theta \in \Theta} P_{\theta'}(x^T)\xi(\theta')}, \quad P_\theta(x^T) = \sum_{c^T \in \mathcal{C}^T} \overbrace{P_\theta(x^T | c^T)}^{\text{Cluster Density}} \underbrace{P_\theta(c^T)}_{\text{Cluster prior}} \quad (2.1)$$

Bayesian formulation of the clustering problem

- ▶ Prior ξ on parameter space Θ .
- ▶ Data $x^T = x_1, \dots, x_T$. Cluster assignments c^T unknown.
- ▶ Posterior $\xi(\cdot | x^T)$.

Posterior distribution

$$\xi(\theta | x^T) = \frac{P_\theta(x^T)\xi(\theta)}{\sum_{\theta \in \Theta} P_{\theta'}(x^T)\xi(\theta')}, \quad P_\theta(x^T) = \sum_{c^T \in \mathcal{C}^T} \overbrace{P_\theta(x^T | c^T)}^{\text{Cluster Density}} \underbrace{P_\theta(c^T)}_{\text{Cluster prior}} \quad (2.1)$$

Marginal posterior prediction

$$P_\xi(c_t | x_t, x^T) = \sum_{\theta \in \Theta} P_\theta(c_t | x_t) \xi(\theta | x^T)$$

Example 4 (Preference clustering)

$$\mathcal{C} = \{1, \dots, C\}, \quad \mathbf{x}_{t,m} \in \{0, 1\}.$$

$$\theta = (\theta_1, \theta_2).$$

Model family

$$P_{\theta_1}(c_t = c) = \theta_{1,c}, \quad c_t \sim \text{Multinomial}(\theta_1) \quad (2.2)$$

$$P_{\theta_2}(x_{t,m} = 1 \mid c_t = c) = \theta_{2,m,c} \quad x_{t,m} \mid c_t = c \sim \text{Bernoulli}(\theta_{2,m,c}) \quad (2.3)$$

Prior

$$\theta_1 \sim \text{Dirichlet}(\gamma), \quad \theta_2 \sim \text{Beta}(\alpha, \beta) \quad (2.4)$$

Supervised learning

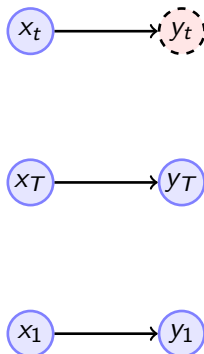


Figure: Graphical model for a classical supervised learning problem.

Semi-supervised learning

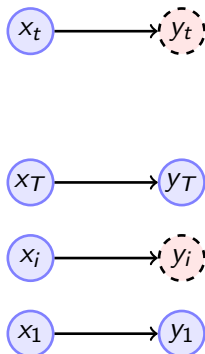


Figure: Graphical model for a classical semi-supervised learning problem.

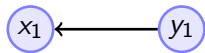
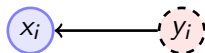
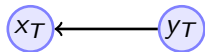
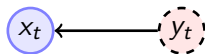


Figure: Generative version of the semi-supervised model

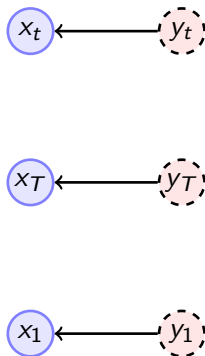


Figure: Basic unsupervised learning model

Applications

- Clustering
- Compression

Recommendation systems

More fun with latent variable models

Social networks

Sequential structures

Network model



Figure: Graphical model for data from a social network.

Network model

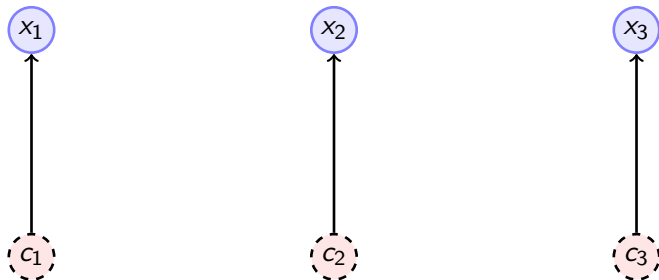


Figure: Graphical model for data from a social network.

Network model

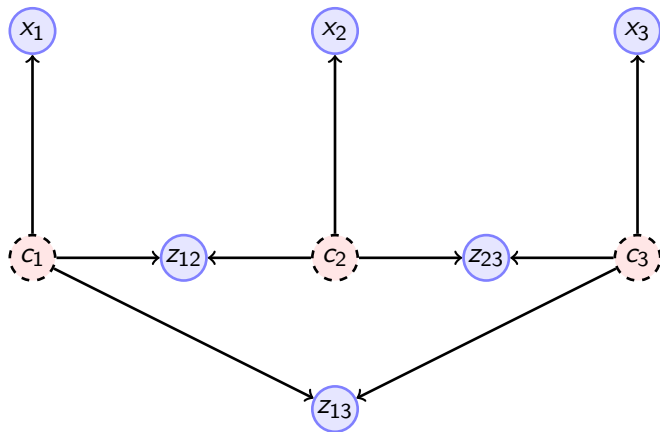
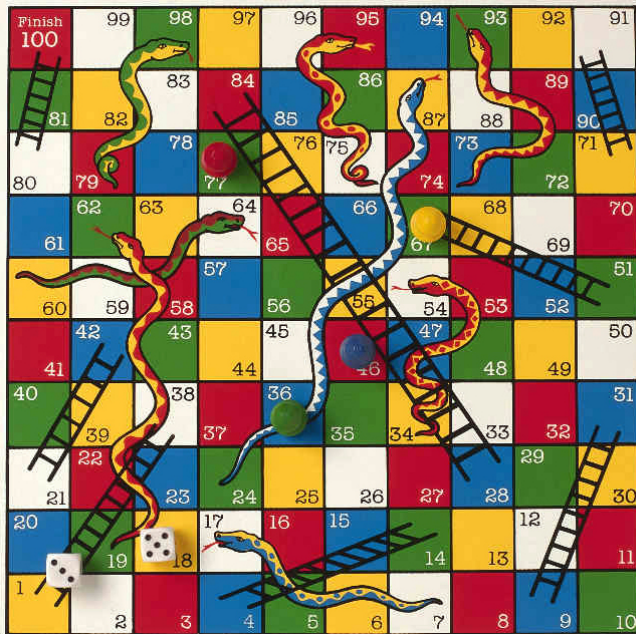


Figure: Graphical model for data from a social network.



Markov process

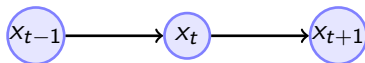


Figure: Graphical model for a Markov process.

Definition 5 (Markov process)

A Markov process is a sequence of variables $x_t : \Omega \rightarrow \mathcal{X}$ such that $x_{t+1} \mid x_t \perp\!\!\!\perp x_{t-k} \forall k \leq 1$.

Application

- ▶ Sequence compression (especially with variable order Models).
- ▶ Web-search (Page-Rank)
- ▶ Hidden Markov Models.
- ▶ MCMC.

Hidden Markov model

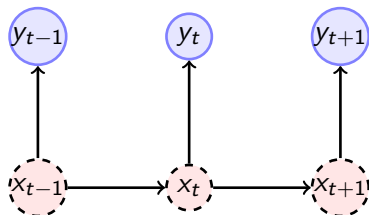


Figure: Graphical model for a hidden Markov model.

$$P_{\theta}(x_{t+1} \mid x_t)$$

(transition distribution)

$$P_{\theta}(y_t \mid x_t)$$

(emission distribution)

Application

- ▶ Speech recognition.
- ▶ Filtering (Kalman Filter).

▶ DNA analysis