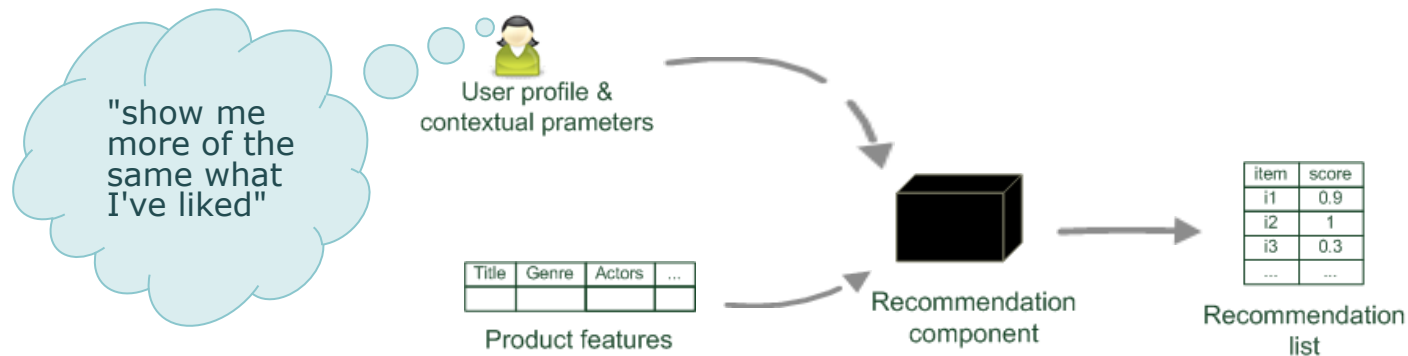# Content-based recommendation

# Content-based recommendation

- **While CF – methods do not require any information about the items,**
  - it might be reasonable to exploit such information; and
  - recommend fantasy novels to people who liked fantasy novels in the past

- **What do we need:**
  - some information about the available items such as the genre ("content")
  - some sort of *user profile* describing what the user likes (the preferences)

- **The task:**
  - learn user preferences
  - locate/recommend items that are "similar" to the user preferences

# What is the "content"?

- **Most CB-recommendation techniques were applied to recommending text documents.**
  - Like web pages or newsgroup messages for example.
  - Now also multimedia content (fashion, music) or e-commerce

- **Content of items can also be represented as text documents.**
  - With textual descriptions of their basic characteristics.
  - Structured: Each item is described by the same set of attributes

| Title | Genre | Author | Type | Price | Keywords |
|---|---|---|---|---|---|
| The Night of the Gun | Memoir | David Carr | Paperback | 29.90 | Press and journalism, drug addiction, personal memoirs, New York |
| The Lace Reader | Fiction, Mystery | Brunonia Barry | Hardcover | 49.90 | American contemporary fiction, detective, historical |
| Into the Fire | Romance, Suspense | Suzanne Brockmann | Hardcover | 45.90 | American fiction, murder, neo-Nazism |

  - Unstructured: free-text description.

# Content representation and item similarities

- **Item representation**

| Title | Genre | Author | Type | Price | Keywords |
|---|---|---|---|---|---|
| The Night of the Gun | Memoir | David Carr | Paperback | 29.90 | Press and journalism, drug addiction, personal memoirs, New York |
| The Lace Reader | Fiction, Mystery | Brunonia Barry | Hardcover | 49.90 | American contemporary fiction, detective, historical |
| Into the Fire | Romance, Suspense | Suzanne Brockmann | Hardcover | 45.90 | American fiction, murder, neo-Nazism |

- **User profile**

| Title | Genre | Author | Type | Price | Keywords |
|---|---|---|---|---|---|
| … | Fiction | Brunonia, Barry, Ken Follett | Paperback | 25.65 | Detective, murder, New York |

$keywords(b_j)$ describes Book $b_j$ with a set of keywords

- **Simple approach**
  - Compute the similarity of an unseen item with the user profile based on the keyword overlap (e.g. using the Dice coefficient)

$$\frac{2 \times \left|keywords(b_i) \cap keywords(b_j)\right|}{\left|keywords(b_i)\right| + \left|keywords(b_j)\right|}$$

  - Or use and combine multiple metrics

# Term-Frequency - Inverse Document Frequency ($TF - IDF$)

- **Simple keyword representation has its problems**
  - in particular when automatically extracted as
    - not every word has similar importance
    - longer documents have a higher chance to have an overlap with the user profile

- **Standard measure: TF-IDF**
  - Encodes text documents in multi-dimensional Euclidian space
    - weighted term vector
  - TF: Measures, how often a term appears (density in a document)
    - assuming that important terms appear more often
    - normalization has to be done in order to take document length into account
  - IDF: Aims to reduce the weight of terms that appear in all documents
    - May not be relevant in some cases (e.g. Male vs. Female attribute on dating sites)
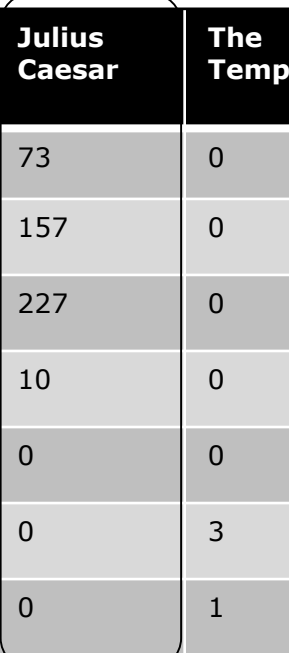
# TF-IDF II

- **Given a keyword $i$ and a document $j$**

- $TF(i, j)$
  - term frequency of keyword $i$ in document $j$

- $IDF(i)$
  - inverse document frequency calculated as $IDF(i) = log \dfrac{N}{n(i)}$
    - $N$ : number of all recommendable documents
    - $n(i)$ : number of documents from $N$ in which keyword $i$ appears

- $TF - IDF$
  - is calculated as: $TF\text{-}IDF(i, j) = TF(i, j) * IDF(i)$

# Example TF-IDF representation

- **Term frequency:**
  - Each document is a count vector in $\mathbb{N}^{|v|}$

| | Antony and Cleopatra | Julius Caesar | The Tempest | Hamlet | Othello | Macbeth |
|---|---|---|---|---|---|---|
| **Antony** | 157 | 73 | 0 | 0 | 0 | 0 |
| **Brutus** | 4 | 157 | 0 | 1 | 0 | 0 |
| **Caesar** | 232 | 227 | 0 | 2 | 1 | 1 |
| **Calpurnia** | 0 | 10 | 0 | 0 | 0 | 0 |
| **Cleopatra** | 57 | 0 | 0 | 0 | 0 | 0 |
| **mercy** | 1.51 | 0 | 3 | 5 | 5 | 1 |
| **worser** | 1.37 | 0 | 1 | 1 | 1 | 0 |

Vector $v$ with dimension $|v| = 7$

Example taken from http://informationretrieval.org

# Example TF-IDF representation

- **Combined TF-IDF weights**
  - Each document is now represented by a real-valued vector of $TF$-$IDF$ weights $\in \mathbb{R}^{|v|}$

| | Antony and Cleopatra | Julius Caesar | The Tempest | Hamlet | Othello | Macbeth |
|---|---|---|---|---|---|---|
| Antony | 157 | 73 | 0 | 0 | 0 | 0 |
| Brutus | 4 | | | | | |
| Caesar | 232 | | | | | |
| Calpurnia | 0 | | | | | |
| Cleopatra | 57 | | | | | |
| mercy | 1.5 | | | | | |
| worser | 1.3 | | | | | |

| | Antony and Cleopatra | Julius Caesar | The Tempest | Hamlet | Othello | Macbeth |
|---|---|---|---|---|---|---|
| Antony | 5.25 | 3.18 | 0 | 0 | 0 | 0.35 |
| Brutus | 1.21 | 6.1 | 0 | 1 | 0 | 0 |
| Caesar | 8.59 | 2.54 | 0 | 1.51 | 0.25 | 0 |
| Calpurnia | 0 | 1.54 | 0 | 0 | 0 | 0 |
| Cleopatra | 2.85 | 0 | 0 | 0 | 0 | 0 |
| mercy | 1.51 | 0 | 1.9 | 0.12 | 5.25 | 0.88 |
| worser | 1.37 | 0 | 0.11 | 4.15 | 0.25 | 1.95 |

Example taken from http://informationretrieval.org

# Improving the vector space model

- **Vectors are usually long and sparse**

- **remove stop words**
  - They will appear in nearly all documents.
  - e.g. "a", "the", "on", …

- **use stemming**
  - Aims to replace variants of words by their common stem
  - e.g. "went" $\Longrightarrow$ "go", "stemming" $\Longrightarrow$ "stem", …

- **size cut-offs**
  - only use top n most representative words to remove "noise" from data
  - e.g. use top 100 words

# Improving the vector space model II

- **Use lexical knowledge, use more elaborate methods for feature selection**
  - Remove words that are not relevant in the domain

- **Detection of phrases as terms**
  - More descriptive for a text than single words
  - e.g. "United Nations"

- **Limitations**
  - semantic meaning remains unknown
  - example: usage of a word in a negative context
    - "there is nothing on the menu that a vegetarian would like.."
    - The word "vegetarian" will receive a higher weight then desired
  - ⇒ an unintended match with a user interested in vegetarian restaurants

# Cosine similarity

- **Usual similarity metric to compare vectors: Cosine similarity (angle)**
  - Cosine similarity is calculated based on the angle between the vectors
    - $sim\left(\vec{a}, \vec{b}\,\right) = \dfrac{\vec{a} \cdot \vec{b}}{|\vec{a}| * |\vec{b}|}$

# Recommending items

- **Simple method: nearest neighbors**

- **May be relevant for item-based recommendations**
  - Most similar items to the currently viewed one

# Query-based retrieval: Rocchio's method

- **Originally for „conversational" query retrieval systems**

- **Query-based retrieval: Rocchio's method**
  - The SMART System: Users are allowed to rate (relevant/irrelevant) retrieved documents (feedback)
  - The system then learns a prototype of relevant/irrelevant documents
  - Queries are then automatically extended with additional terms/weight of relevant documents

- **The paradigm fits well also for recommender systems**

# Rocchio details

- **Document collections $D^+$ (liked) and $D^-$ (disliked)**
  - Calculate prototype vector for these categories.



| | |
|---|---|
| ○ | Relevant documents |
| X | Nonrelevant documents |
| ◑ | Centroids |

- **Computing modified query $Q_{i+1}$ from current query $Q_i$ with:**

$$Q_{i+1} = \alpha * Q_i + \beta \left( \frac{1}{|D^+|} \sum_{d^+ \in D^+} d^+ \right) - \gamma \left( \frac{1}{|D^-|} \sum_{d^- \in D^-} d^- \right)$$

- **$\alpha, \beta, \gamma$ used to fine-tune the feedback**
  - $\alpha$ weight for original query
  - $\beta$ weight for positive feedback
  - $\gamma$ weight for negative feedback

- **Often only positive feedback is used**
  - More valuable than negative feedback

# Practical challenges of Rocchio's method

- **Certain number of item ratings needed to build reasonable user model**
  - Can be automated by trying to capture user ratings implicitly (click on document)
  - Pseudorelevance Feedback: Assume that the first $n$ documents match the query best. The set $D^-$ is not used until explicit negative feedback exists.

- **User interaction required during retrieval phase**
  - Interactive query refinement opens new opportunities for gathering information and
  - Helps user to learn which vocabulary should be used to receive the information he needs

# Explicit decision models

- **Decision tree for recommendation problems**
  - inner nodes labeled with item features (keywords)
  - used to partition the test examples
    - **existence or non existence of a keyword**
  - in basic setting only two classes appear at leaf nodes
    - **interesting or not interesting**
  - decision tree can automatically be constructed from training data
  - works best with small number of features
  - use meta features like author name, genre, ...  instead of TF-IDF representation.

# Explicit decision models II

- **Rule induction**
  - built on RIPPER algorithm
  - good performance compared with other classification methods
    - **eloborate postpruning techniques of RIPPER**
    - **extension for e-mail classification**
      - **takes document structure into account**

- **main advantages of these decision models:**
  - inferred decision rules serve as basis for generating explanations for recommendation
  - existing domain knowledge can be incorporated in models

# On feature selection

- **process of choosing a subset of available terms**

- **different strategies exist for deciding which features to use**
  - feature selection based on domain knowledge and lexical information from WordNet (Pazzani and Billsus 1997)
  - frequency-based feature selection to remove words appearing "too rare" or "too often" (Chakrabarti 2002)

- **Not appropriate for larger text corpora**
  - Better to
    - evaluate value of individual features (keywords) independently and
    - construct a ranked list of "good" keywords.

- **Typical measure for determining utility of keywords: e.g. $X^2$, mutual information measure or Fisher's discrimination index**

# Limitations of content-based recommendation methods

- **Keywords alone may not be sufficient to judge quality/relevance of a document or web page**
    - up-to-date-ness, usability, aesthetics, writing style
    - content may also be limited / too short
    - content may not be automatically extractable (multimedia)

- **Ramp-up phase required**
    - Some training data is still required
    - Web 2.0: Use other sources to learn the user preferences

- **Overspecialization**
    - Algorithms tend to propose "more of the same"
    - Or: too similar news items
    - Multicriterial optimization (diversity, novelty)

# Discussion & summary

- In contrast to collaborative approaches, content-based techniques do not require user community in order to work

- Presented approaches aim to learn a model of user's interest preferences based on explicit or implicit feedback
  - Deriving implicit feedback from user behavior can be problematic

- Evaluations show that a good recommendation accuracy can be achieved with help of machine learning techniques
  - These techniques do not require a user community

- Danger exists that recommendation lists contain too many similar items
  - All learning techniques require a certain amount of training data
  - Some learning methods tend to overfit the training data

- Pure content-based systems are rarely found in commercial environments

# Knowledge-based recommendation

# Basic I/O Relationship

**Knowledge-based: "Tell me what fits based on my needs"**

User profile & contextual prameters

| Title | Genre | Actors | ... |
|---|---|---|---|
|  |  |  |  |

Product features

Knowledge models

Recommendation component

| item | score |
|---|---|
| i1 | 0.9 |
| i2 | 1 |
| i3 | 0.3 |
| ... | ... |

Recommendation list

# Why do we need knowledge-based recommendation?

- **Products with low number of available ratings**



- **Time span plays an important role**
  - five-year-old ratings for computers
  - user lifestyle or family situation changes

- **Customers want to define their requirements explicitly**
  - "the color of the car should be black"

# Knowledge-based recommender systems

- **Constraint-based**
  - based on explicitly defined set of recommendation rules
  - *(partially)* fulfill recommendation rules

- **Case-based**
  - Item-based: give me similar items, however with larger display

- **Both approaches are similar in their conversational recommendation proces** *(edge of query retrieval and recommender systems)*
  - users specify the requirements
  - systems try to identify solutions
  - if no solution can be found, users change requirements
  - *Not always, we may learn knowledge RS rules from collaborative data*

# Constraint-based recommender systems

- **Knowledge base**
  - usually mediates between user model and item properties
  - variables
    - user model features (requirements), Item features (catalogue)
  - set of constraints
    - logical implications (IF user requires A THEN proposed item should possess feature B)
    - hard and soft/weighted constraints
    - solution preferences

- **Derive a set of recommendable items**
  - fulfilling set of applicable constraints
  - applicability of constraints depends on current user model
  - explanations – transparent line of reasoning

# Constraint-based recommendation tasks

- **Find a set of user requirements such that a subset of items fulfills all constraints**
  - ask user which requirements should be relaxed/modified such that some items exist that do not violate any constraint

- **Find a subset of items that satisfy the maximum set of weighted constraints**
  - similar to find a maximally succeeding subquery (XSS)
  - all proposed items have to fulfill the same set of constraints
  - compute relaxations based on predetermined weights

- **Rank items according to weights of satisfied soft constraints**
  - rank items based on the ratio of fulfilled constraints
  - does not require additional ranking scheme

# Ranking the items

- **Multi-attribute utility theory**
  - each item is evaluated according to a predefined set of dimensions that provide an aggregated view on the basic item properties

- *E.g. quality and economy are dimensions in* **the domain of digital cameras**

| id | value | quality | economy |
|----|-------|---------|---------|
| price | ≤250 | 5 | 10 |
|  | >250 | 10 | 5 |
| mpix | ≤8 | 4 | 10 |
|  | >8 | 10 | 6 |
| opt-zoom | ≤9 | 6 | 9 |
|  | >9 | 10 | 6 |
| LCD-size | ≤2.7 | 6 | 10 |
|  | >2.7 | 9 | 5 |
| movies | Yes | 10 | 7 |
|  | no | 3 | 10 |
| sound | Yes | 10 | 8 |
|  | no | 7 | 10 |
| waterproof | Yes | 10 | 6 |
|  | no | 8 | 10 |

# Item utility for customers

- **Customer specific interest**

| Customer | quality | economy |
|---|---|---|
| $Cu_1$ | 80% | 20% |
| $Cu_2$ | 40% | 60% |

- ***Calculation of Utility***

| quality | economy | $cu_1$ | $cu_2$ |
|---|---|---|---|
| $P_1$ Σ(5,4,6,6,3,7,10) = 41 | Σ (10,10,9,10,10,10,6) = 65 | 45.8 [8] | 55.4 [6] |
| $P_2$ Σ(5,4,6,6,10,10,8) = 49 | Σ (10,10,9,10,7,8,10) = 64 | 52.0 [7] | 58.0 [1] |
| $P_3$ Σ(5,4,10,6,10,10,8) = 53 | Σ (10,10,6,10,7,8,10) = 61 | 54.6 [5] | 57.8 [2] |
| $P_4$ Σ(5,10,10,6,10,7,10) = 58 | Σ (10,6,6,10,7,10,6) = 55 | 57.4 [4] | 56.2 [4] |
| $P_5$ Σ(5,4,6,10,10,10,8) = 53 | Σ (10,10,9,6,7,8,10) = 60 | 54.4 [6] | 57.2 [3] |
| $P_6$ Σ(5,10,6,9,10,10,8) = 58 | Σ (10,6,9,5,7,8,10) = 55 | 57.4 [3] | 56.2 [5] |
| $P_7$ Σ(10,10,6,9,10,10,8) = 63 | Σ (5,6,9,5,7,8,10) = 50 | 60.4 [2] | 55.2 [7] |
| $P_8$ Σ(10,10,10,9,10,10,10) = 69 | Σ (5,6,6,5,7,8,6) = 43 | 63.8 [1] | 53.4 [8] |

# Case-based recommender systems

- **Items are retrieved using similarity measures**

- **Distance similarity**

$$similarity(p, REQ) = \frac{\sum_{r \in REQ} w_r * sim(p, r)}{\sum_{r \in REQ} w_r}$$

- **Def.**
  - sim (p, r) expresses for each item attribute value φr (p) its distance to the customer requirement r ∈ REQ.
  - $w_r$ is the importance weight for requirement r

- **In real world, customer  would like to**
  - maximize certain properties. i.e. resolution of a camera, "more is better"(MIB)
  - minimize certain properties. i.e. price of a camera, "less is better"(LIB)
  - Target within some values, e.g. Price between x,y

# Constraint-based recommendation tasks

# Knowledge-based recommender systems

- **Transform known rating on items into**
  - Rating (preference regression) of item features
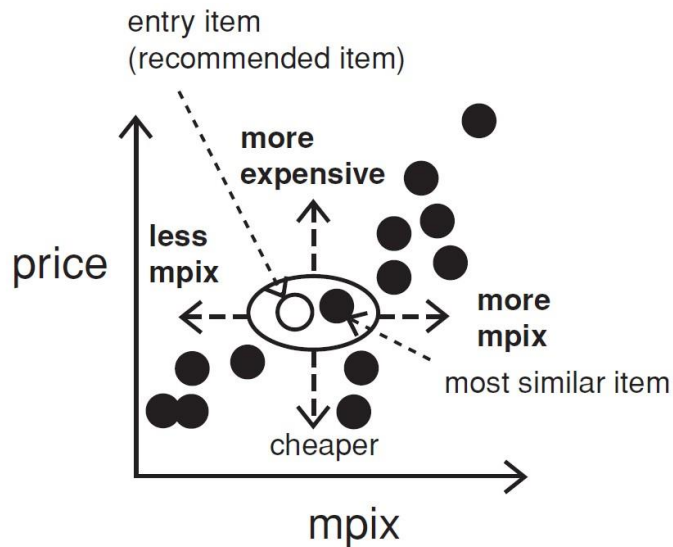


  - Learning combination of item feature´s ratings
    - Based on goodness of fit on features

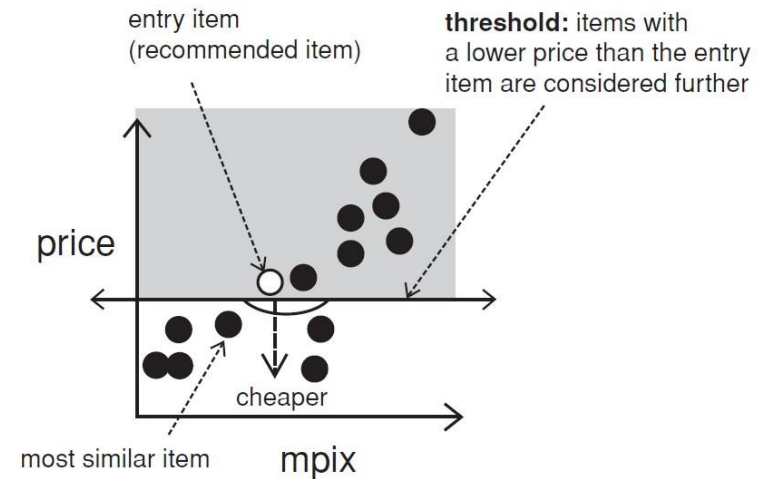$$@(o) = \frac{(2 * f_{Price}(o) + 1 * f_{Display}(o) + 1 * f_{RAM}(o))}{4}$$

- **Evaluate the learned rating function on all other objects**
  - *Recommend better instead of similar objects*

# Interacting with case-based recommenders

- **Customers maybe not know what they are seeking**

- **Critiquing is an effective way to support such navigations**

- **Customers specify their change requests (*price or mpix*) that are not satisfied by the current item (*entry item*)**
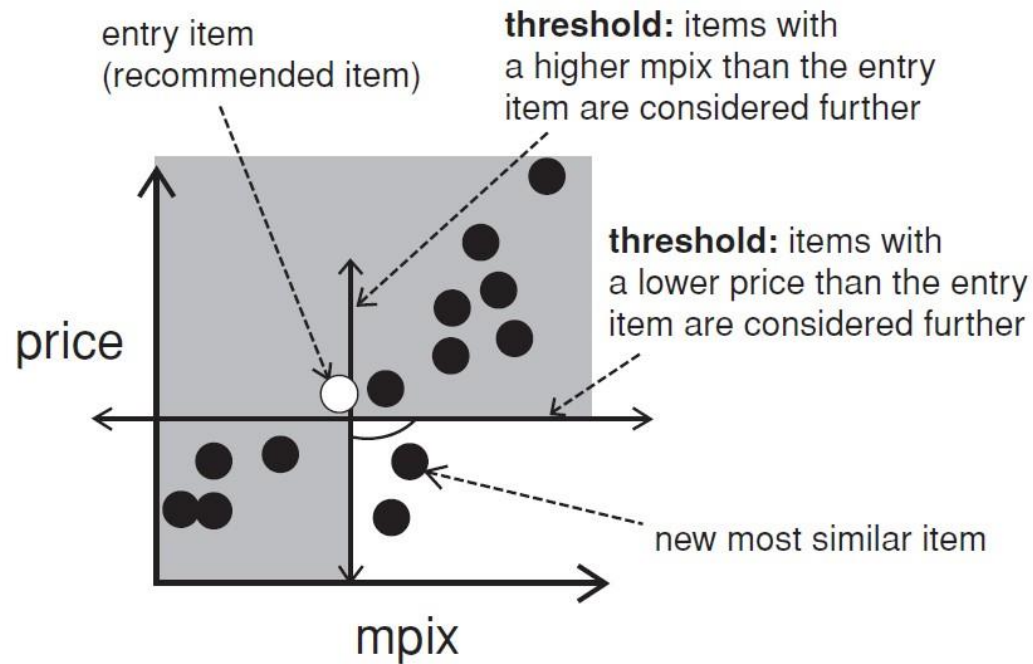


*Critique on price*

# Compound critiques

- **Operate over multiple properties can improve the efficiency of recommendation dialogs**

# Summary

- **Knowledge-based recommender systems**
  - Move from recommending *similar* to recommending *better* objects

- **Limitations**
  - cost of knowledge acquisition
    - from domain experts
    - from users
    - from web resources
  - accuracy of preference models
    - very fine granular preference models require many interaction cycles
    - collaborative filtering models preference implicitly
  - independence assumption can be challenged
    - **preferences are not always independent from each other**