

Matrix Factorization

Used for:

- dim. red
- imputation of missing data
- denoising by "borrowing strength"
- finding "archetypes" or "latent dims"

- looking to find a low rank matrix

st: loss is minimized (for real valued

matrices, sq. loss typically used). Using sq. loss is
same as minimizing Frob. norm.

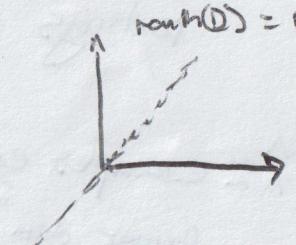
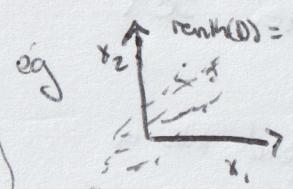
e.g. PCA

$$\hat{D} = \underset{\hat{D} \in \mathbb{R}^{m \times n}}{\operatorname{argmin}} \| D - \hat{D} \|_F^2$$

$$\operatorname{rank}(\hat{D}) = k$$

where $D \in \mathbb{R}^{m \times n}$ is the design matrix

$\operatorname{rank}(D) = \dim. \text{ of subspace spanned by rows in } D$



Rank constraint can be implicitly encoded by letting $\tilde{D} = XY$

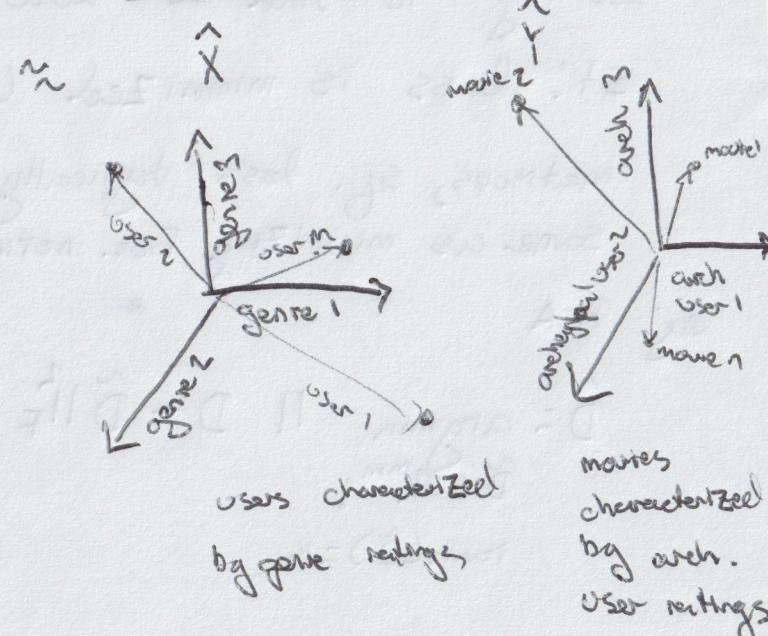
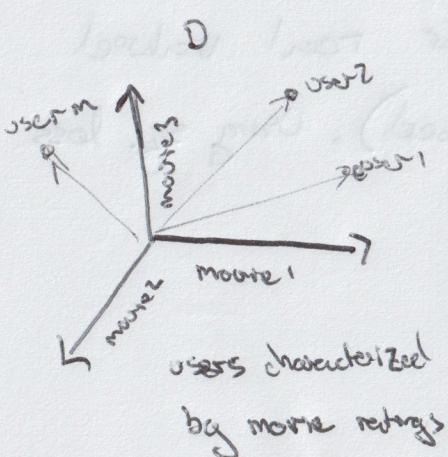
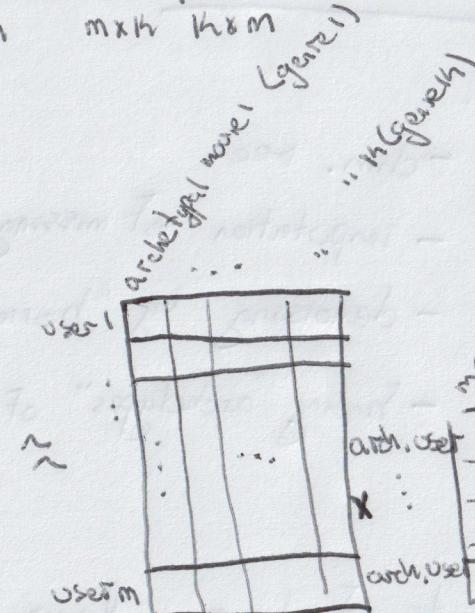
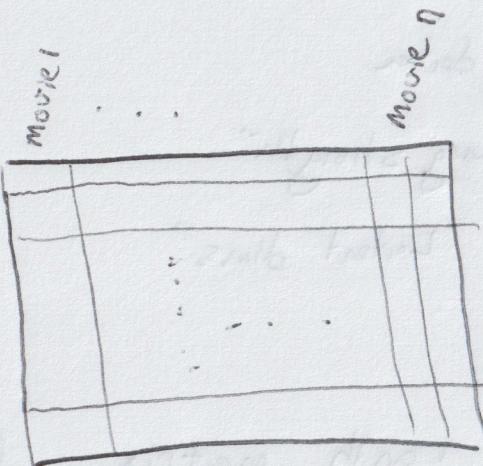
$$\Rightarrow \hat{X}, \hat{Y} = \underset{\substack{X \in \mathbb{R}^{m \times k} \\ Y \in \mathbb{R}^{k \times n}}}{\text{arg min}} \|D - XY\|_F^2$$

$$\hat{D} = \hat{X}\hat{Y}$$

e.g. a user \times movie ratings dataset

\hat{X} ~ examples embedded into a feature space

\hat{Y} ~ features " lower dim. example space



in general, we can add regularization to keep generalization error low as well as summing over only observed data in a sparse dataset, as well as add desired constraints (e.g. non-negativity of \hat{D})

$$\hat{X}, \hat{Y} = \underset{\substack{X \in \mathbb{R}^{m \times n} \\ Y \in \mathbb{R}^{n \times k} \\ \text{constraints}}}{\operatorname{arg\,min}} \sum_{(i,j) \in \Omega} (D_{ij} - X_i^T Y_j)^2 + \lambda \sum_{i=1}^m \Gamma_X(X_i) + \lambda \sum_{j=1}^k \Gamma_Y(Y_j)$$

$\left. \begin{array}{c} \text{rows of } X \\ \text{cols of } Y \end{array} \right\}$

 $\Gamma_X : \mathbb{R}^n \rightarrow \mathbb{R}_+$
 (reg. on rows of X)

 $\Gamma_Y : \mathbb{R}^k \rightarrow \mathbb{R}_+$
 (reg. on cols of Y)

$\Omega = \{(i,j) \text{ indices of } D \text{ which are observed}\}$

$$\hat{D}_{ij} = \hat{X}_i^T \hat{Y}_j$$

- the optimization is usually solved w/

Alternating minimization (hold X constant and optimize over Y , hold Y constant and optimize over X ...), similar to coordinate descent

OR

SGD (each iteration takes $O(\text{batch size})$)

vs.

(each iteration takes $O(|\Omega|)$ for Alt. min.)

- general recipe for Matrix Fact.
- choose loss function
- choose regularization
- impose any constraints
- choose optimization algo.

Generalized Low Rank Models

e.g.: booleans, ordinals, categoricals

- is applicable to design matrices w/ any data type (not just reals)
- can embed the data into \mathbb{R}^n by finding the "best" rank k
real valued matrix
- can also project^{matrix} back into original data types if desired

let $F_j = \text{set of values that col. } j \text{ of } D \text{ can take on}$ (e.g.: $F_j = \{-1, 1\}$ for booleans)

and let $L_j: F_j \times \mathbb{R} \rightarrow \mathbb{R}_+$ be an appropriate loss function for that data type. Then under a GLM

$$\hat{X}, \hat{Y} = \underset{\substack{X \in \mathbb{R}^{n \times k} \\ Y \in \mathbb{R}^{k \times n} \\ \text{constraints}}}{\operatorname{argmin}} \sum_{(i,j) \in \Omega} L_j(D_{ij}, X_i^T y_j) + \lambda \sum_{i=1}^m \Gamma_x(x_i) + \gamma \sum_{i=1}^n \Gamma_y(y_i)$$

- then \hat{X}, \hat{Y} is the "best" low rank representation of the data in \mathbb{R}^n . \hat{X} is the "best" representation of the examples in \mathbb{R}^k

- if desired, to express the dexter in the original dexter types

$$\hat{D}_{ij} = \underset{d \in \mathcal{F}_j}{\operatorname{argmin}} L_j(d, \hat{x}^T, \hat{g}_j)$$

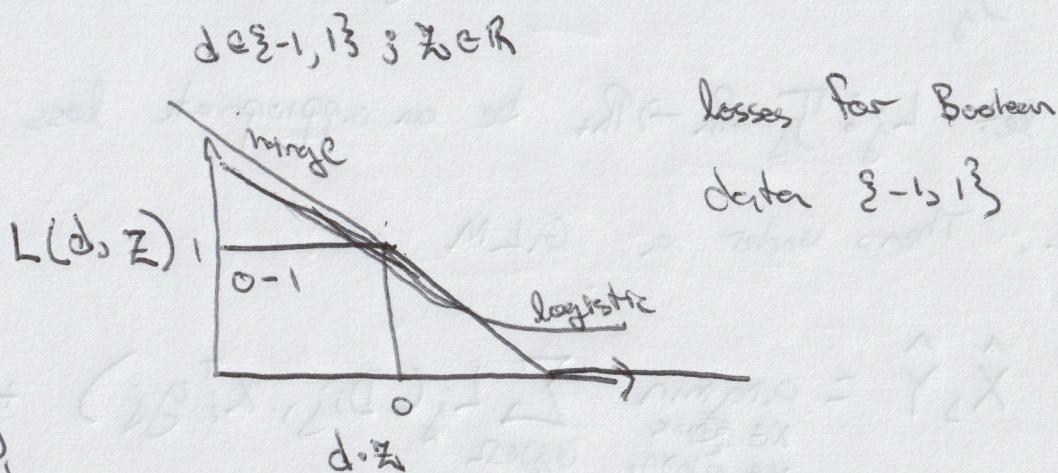
- that is, we choose the value of \hat{D}_{ij} s.t. the loss between this choice and the real-valued, low rank approx' of the original design matrix is minimized

Possible loss functions

Boolean: hinge, logistic

categorical: log loss / cross entropy

ordinal: encode levels to numbers and use ordinal hinge-loss



$$L_{\text{hinge}}(d, z) = \max(0, 1 - dz)$$

$$L_{\text{logistic}}(d, z) = \log(1 + e^{-dz})$$