

Collaborative filtering in Hilbert spaces with spectral regularization

Jacob Abernethy¹ Francis Bach²
Theodoros Evgeniou³ **Jean-Philippe Vert⁴**

¹UC Berkeley

²INRIA / Ecole normale superieure de Paris

³INSEAD

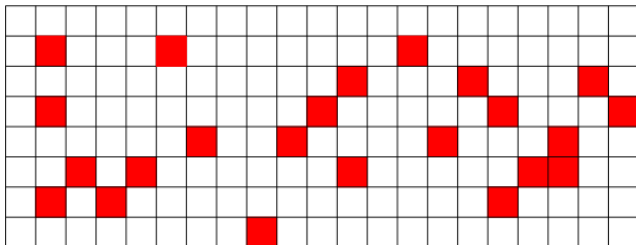
⁴ParisTech / Institut Curie / INSERM

Second Canada-France Congress of Mathematics, Montreal,
Canada, June 3, 2008.

Collaborative Filtering (CF)

The problem

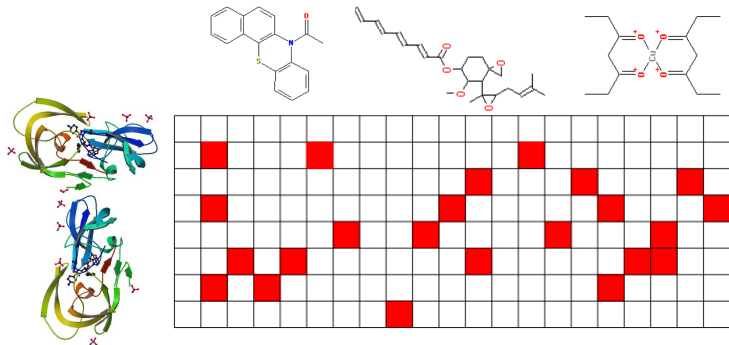
- Given a set of n_x “movies” $\mathbf{x} \in \mathcal{X}$ and a set of n_y “customers” $\mathbf{y} \in \mathcal{Y}$,
- predict the “rating” $z(\mathbf{x}, \mathbf{y}) \in \mathcal{Z}$ of customer \mathbf{y} for movie \mathbf{x}
- Training data: large $n_x \times n_y$ incomplete matrix Z that describes the known ratings of some customers for some movies
- Goal: complete the matrix.



Another CF example

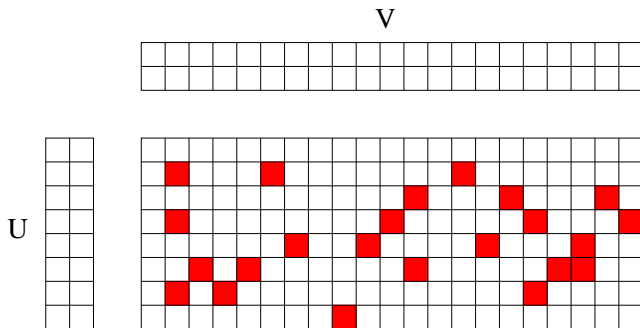
Drug design

- Given a family of **proteins** of therapeutic interest (e.g., GPCR's)
- Given all known **small molecules** that bind to these proteins
- Can we predict unknown **interactions**?



CF by low-rank matrix approximation

- A common strategy for CF
- Z has rank less than $k \Leftrightarrow Z = UV^T$ $U \in \mathbb{R}^{n_x \times k}$, $V \in \mathbb{R}^{n_y \times k}$
- Examples: PLSA (Hoffmann, 2001), MMMF (Srebro et al, 2004)
- Numerical and statistical efficiency



CF by low-rank matrix approximation example

Fitting low-rank models (Srebro et al, 2004)

- **Relax** the (non-convex) rank of Z into the (**convex**) **trace norm** of Z : if $\sigma_i(Z)$ are the singular values of Z ,

$$\text{rank}Z = \sum_i \mathbf{1}_{\sigma_i(Z)>0} \qquad \|Z\|_* = \sum_i \sigma_i(Z).$$

- n observations z_u corresponding to $\mathbf{x}_{i(u)}$ and $\mathbf{y}_{j(u)}$, $u = 1, \dots, n$:

$$\min_{Z \in \mathbb{R}^{n \times n_y}} \sum_{u=1}^n \ell(z_u, Z_{i(u), j(u)}) + \lambda \|Z\|_*,$$

where $\ell(z, z')$ is a convex loss function.

- This is an SDP if ℓ is SDP-representable

Fitting low-rank models (Srebro et al, 2004)

- **Relax** the (non-convex) rank of Z into the (**convex**) **trace norm** of Z : if $\sigma_i(Z)$ are the singular values of Z ,

$$\text{rank}Z = \sum_i \mathbf{1}_{\sigma_i(Z)>0} \qquad \|Z\|_* = \sum_i \sigma_i(Z).$$

- n observations z_u corresponding to $\mathbf{x}_{i(u)}$ and $\mathbf{y}_{j(u)}$, $u = 1, \dots, n$:

$$\min_{Z \in \mathbb{R}^{n_x \times n_y}} \sum_{u=1}^n \ell(z_u, Z_{i(u), j(u)}) + \lambda \|Z\|_*,$$

where $\ell(z, z')$ is a convex loss function.

- This is an SDP if ℓ is SDP-representable

The problem

- Often we have **additional attributes**:
 - gender, age of customers; type, actors of movies..
 - 3D structures of proteins and ligands for protein-ligand interaction prediction
- **How to include attributes in CF?**
- Expected gains: increase **performance**, allow predictions on **new** movie and/or customers.

Our contributions

- A **general framework** for CF **with or without attributes**, using **kernels** to describe attributes (“kernel-CF”)
- A **family of algorithms** for CF in this setting

The problem

- Often we have **additional attributes**:
 - gender, age of customers; type, actors of movies..
 - 3D structures of proteins and ligands for protein-ligand interaction prediction
- **How to include attributes in CF?**
- Expected gains: increase **performance**, allow predictions on **new** movie and/or customers.

Our contributions

- A **general framework** for CF **with or without attributes**, using **kernels** to describe attributes (“kernel-CF”)
- A **family of algorithms** for CF in this setting

Basic facts

- n_x movies and n_y customers
- The known rating $z(\mathbf{x}_i, \mathbf{y}_j)$ of customer \mathbf{y}_j for movie \mathbf{x}_i is stored in the (i, j) -th entry of a **matrix** M (of size $n_x \times n_y$).
- M represents a **linear application** / **bilinear form**:

$$M : \mathbb{R}^{n_y} \rightarrow \mathbb{R}^{n_x}$$

defined by:

$$\mathbf{e}_i^\top M \mathbf{f}_j = M_{i,j}$$

- Rank / trace norm are **spectral properties** of the linear application

The idea

Reformulations

- **Represent** the i -th movie $\mathbf{x}_i \in \mathcal{X}$ (resp. j -th customer $\mathbf{y}_j \in \mathcal{Y}$) by the i -th basis vector $\mathbf{e}_i \in \mathbb{R}^{n_x}$ (resp. $\mathbf{f}_j \in \mathbb{R}^{n_y}$):

$$\phi_X(\mathbf{x}_i) = \mathbf{e}_i, \quad \phi_Y(\mathbf{y}_j) = \mathbf{f}_j.$$

- **Approximate** the rating function by a **bilinear form**:

$$\forall (\mathbf{x}_i, \mathbf{y}_j) \in \mathcal{X} \times \mathcal{Y}, \quad G_M(\mathbf{x}_i, \mathbf{y}_j) = \phi_X(\mathbf{x}_i)^\top M \phi_Y(\mathbf{y}_j),$$

by constraining a **spectral property** of $M : \mathbb{R}^{n_x} \mapsto \mathbb{R}^{n_x}$.

An idea

If we have additional attributes about movies / customer, why not include them in $\phi_X(\mathbf{x})$ and $\phi_Y(\mathbf{y})$?

The idea

Reformulations

- **Represent** the i -th movie $\mathbf{x}_i \in \mathcal{X}$ (resp. j -th customer $\mathbf{y}_j \in \mathcal{Y}$) by the i -th basis vector $\mathbf{e}_i \in \mathbb{R}^{n_x}$ (resp. $\mathbf{f}_j \in \mathbb{R}^{n_y}$):

$$\phi_X(\mathbf{x}_i) = \mathbf{e}_i, \quad \phi_Y(\mathbf{y}_j) = \mathbf{f}_j.$$

- **Approximate** the rating function by a **bilinear form**:

$$\forall (\mathbf{x}_i, \mathbf{y}_j) \in \mathcal{X} \times \mathcal{Y}, \quad G_M(\mathbf{x}_i, \mathbf{y}_j) = \phi_X(\mathbf{x}_i)^\top M \phi_Y(\mathbf{y}_j),$$

by constraining a **spectral property** of $M : \mathbb{R}^{n_x} \mapsto \mathbb{R}^{n_x}$.

An idea

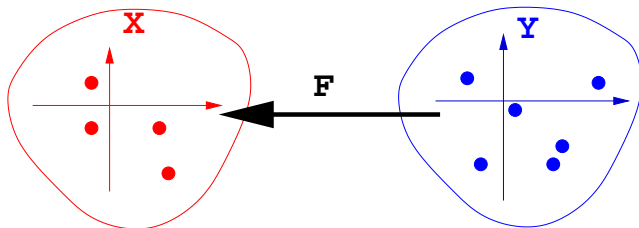
If we have additional attributes about movies / customer, why not include them in $\phi_X(\mathbf{x})$ and $\phi_Y(\mathbf{y})$?

Setting

- Movies: points in a Hilbert space \mathcal{X}
- Customers: points in a Hilbert space \mathcal{Y}
- We model the preference of customer \mathbf{y} for a movie \mathbf{x} by a bilinear form:

$$f(\mathbf{x}, \mathbf{y}) = \langle \mathbf{x}, F\mathbf{y} \rangle_{\mathcal{X}},$$

where $F \in \mathcal{B}_0(\mathcal{Y}, \mathcal{X})$ is a **compact linear operator** (i.e., a “matrix”).



Classical results

- For (\mathbf{x}, \mathbf{y}) in $\mathcal{X} \times \mathcal{Y}$ the **tensor product** $\mathbf{x} \otimes \mathbf{y}$ is the operator

$$\forall \mathbf{h} \in \mathcal{Y}, \quad (\mathbf{x} \otimes \mathbf{y}) \mathbf{h} = \langle \mathbf{y}, \mathbf{h} \rangle_{\mathcal{Y}} \mathbf{x}.$$

- Any compact operator $F : \mathcal{Y} \rightarrow \mathcal{X}$ admits a spectral decomposition:

$$F = \sum_{i=1}^{\infty} \sigma_i \mathbf{u}_i \otimes \mathbf{v}_i.$$

where the $\sigma_i \geq 0$ are the **singular values** and $(\mathbf{u}_i)_{i \in \mathbb{N}}$ and $(\mathbf{v}_i)_{i \in \mathbb{N}}$ are orthonormal families in \mathcal{X} and \mathcal{Y} .

- The **spectrum of** F is the set of singular values sorted in decreasing order: $\sigma_1(F) \geq \sigma_2(F) \geq \dots \geq 0$.
- This is the natural generalization of singular values for matrices.

Classical results

- For (\mathbf{x}, \mathbf{y}) in $\mathcal{X} \times \mathcal{Y}$ the **tensor product** $\mathbf{x} \otimes \mathbf{y}$ is the operator

$$\forall \mathbf{h} \in \mathcal{Y}, \quad (\mathbf{x} \otimes \mathbf{y}) \mathbf{h} = \langle \mathbf{y}, \mathbf{h} \rangle_{\mathcal{Y}} \mathbf{x}.$$

- Any compact operator $F : \mathcal{Y} \rightarrow \mathcal{X}$ admits a spectral decomposition:

$$F = \sum_{i=1}^{\infty} \sigma_i \mathbf{u}_i \otimes \mathbf{v}_i.$$

where the $\sigma_i \geq 0$ are the **singular values** and $(\mathbf{u}_i)_{i \in \mathbb{N}}$ and $(\mathbf{v}_i)_{i \in \mathbb{N}}$ are orthonormal families in \mathcal{X} and \mathcal{Y} .

- The **spectrum of** F is the set of singular values sorted in decreasing order: $\sigma_1(F) \geq \sigma_2(F) \geq \dots \geq 0$.
- This is the natural generalization of singular values for matrices.

Classical results

- For (\mathbf{x}, \mathbf{y}) in $\mathcal{X} \times \mathcal{Y}$ the **tensor product** $\mathbf{x} \otimes \mathbf{y}$ is the operator

$$\forall \mathbf{h} \in \mathcal{Y}, \quad (\mathbf{x} \otimes \mathbf{y}) \mathbf{h} = \langle \mathbf{y}, \mathbf{h} \rangle_{\mathcal{Y}} \mathbf{x}.$$

- Any compact operator $F : \mathcal{Y} \rightarrow \mathcal{X}$ admits a spectral decomposition:

$$F = \sum_{i=1}^{\infty} \sigma_i \mathbf{u}_i \otimes \mathbf{v}_i.$$

where the $\sigma_i \geq 0$ are the **singular values** and $(\mathbf{u}_i)_{i \in \mathbb{N}}$ and $(\mathbf{v}_i)_{i \in \mathbb{N}}$ are orthonormal families in \mathcal{X} and \mathcal{Y} .

- The **spectrum of** F is the set of singular values sorted in decreasing order: $\sigma_1(F) \geq \sigma_2(F) \geq \dots \geq 0$.
- This is the natural generalization of singular values for matrices.

Useful classes for operators

Operators of finite rank

- The **rank** of an operator is the number of strictly positive singular values.
- Hence operators of rank smaller or equal to k are characterized by:

$$\sigma_{k+1}(F) = 0.$$

Trace-class operators

The **trace-class** operators are the compact operators F that satisfy:

$$\|F\|_* := \sum_{i=1}^{\infty} \sigma_i(F) < \infty.$$

$\|F\|_*$ is a norm over the trace-class operators, called the **trace norm**.

Operators of finite rank

- The **rank** of an operator is the number of strictly positive singular values.
- Hence operators of rank smaller or equal to k are characterized by:

$$\sigma_{k+1}(F) = 0.$$

Trace-class operators

The **trace-class** operators are the compact operators F that satisfy:

$$\|F\|_* := \sum_{i=1}^{\infty} \sigma_i(F) < \infty.$$

$\|F\|_*$ is a norm over the trace-class operators, called the **trace norm**.

Hilbert-Schmidt operators

- The **Hilbert-Schmidt operators** are compact operators F that satisfy:

$$\|F\|_{Fro}^2 := \sum_{i=1}^{\infty} \sigma_i(F)^2 < \infty.$$

- They form a **Hilbert space** with inner product:

$$\langle \mathbf{x} \otimes \mathbf{y}, \mathbf{x}' \otimes \mathbf{y}' \rangle_{\mathcal{X} \otimes \mathcal{Y}} = \langle \mathbf{x}, \mathbf{x}' \rangle_{\mathcal{X}} \langle \mathbf{y}, \mathbf{y}' \rangle_{\mathcal{Y}}.$$

Definition

A function $\Omega : \mathcal{B}_0(\mathcal{Y}, \mathcal{X}) \mapsto \mathbb{R} \cup \{+\infty\}$ is called a **spectral penalty function** if it can be written as:

$$\Omega(F) = \sum_{i=1}^{\infty} s_i(\sigma_i(F)) ,$$

where for any $i \geq 1$, $s_i : \mathbb{R}^+ \mapsto \mathbb{R}^+ \cup \{+\infty\}$ is a **non-decreasing** penalty function satisfying **$s_i(0) = 0$** .

Spectral penalty function

Examples

- **Rank constraint:** take $s_{k+1}(0) = 0$ and $s_{k+1}(u) = +\infty$ for $u > 0$, and $s_i = 0$ for $i \geq k$. Then

$$\Omega(F) = \begin{cases} 0 & \text{if } \text{rank}(F) \leq k, \\ +\infty & \text{if } \text{rank}(F) > k. \end{cases}$$

- **Trace norm:** take $s_i(u) = u$ for all i , then:

$$\Omega(F) = \|F\|_*.$$

- **Hilbert-Schmidt norm:** take $s_i(u) = u^2$ for all i , then

$$\Omega(F) = \|F\|_{Fro}^2.$$

Examples

- **Rank constraint:** take $s_{k+1}(0) = 0$ and $s_{k+1}(u) = +\infty$ for $u > 0$, and $s_i = 0$ for $i \geq k$. Then

$$\Omega(F) = \begin{cases} 0 & \text{if } \text{rank}(F) \leq k, \\ +\infty & \text{if } \text{rank}(F) > k. \end{cases}$$

- **Trace norm:** take $s_i(u) = u$ for all i , then:

$$\Omega(F) = \|F\|_*.$$

- **Hilbert-Schmidt norm:** take $s_i(u) = u^2$ for all i , then

$$\Omega(F) = \|F\|_{Fro}^2.$$

Examples

- **Rank constraint:** take $s_{k+1}(0) = 0$ and $s_{k+1}(u) = +\infty$ for $u > 0$, and $s_i = 0$ for $i \geq k$. Then

$$\Omega(F) = \begin{cases} 0 & \text{if } \text{rank}(F) \leq k, \\ +\infty & \text{if } \text{rank}(F) > k. \end{cases}$$

- **Trace norm:** take $s_i(u) = u$ for all i , then:

$$\Omega(F) = \|F\|_*.$$

- **Hilbert-Schmidt norm:** take $s_i(u) = u^2$ for all i , then

$$\Omega(F) = \|F\|_{Fro}^2.$$

Learning operator with spectral regularization

Setting

- **Training set:** $(\mathbf{x}_i, \mathbf{y}_i, t_i)_{i=1, \dots, N}$ a set of (movie, customer, preference).
- **Loss function** $l(t, t')$: cost of predicting preference t instead of t' .
- **Empirical risk** of an operator F :

$$R_N(F) = \frac{1}{N} \sum_{i=1}^N l(\langle \mathbf{x}_i, F \mathbf{y}_i \rangle_{\mathcal{X}}, t_i) .$$

Learning an operator

$$\min_{F \in \mathcal{B}_0(\mathcal{Y}, \mathcal{X}), \Omega(F) < \infty} \{ R_N(F) + \lambda \Omega(F) \} .$$

Learning operator with spectral regularization

Setting

- **Training set:** $(\mathbf{x}_i, \mathbf{y}_i, t_i)_{i=1, \dots, N}$ a set of (movie, customer, preference).
- **Loss function** $l(t, t')$: cost of predicting preference t instead of t' .
- **Empirical risk** of an operator F :

$$R_N(F) = \frac{1}{N} \sum_{i=1}^N l(\langle \mathbf{x}_i, F\mathbf{y}_i \rangle_{\mathcal{X}}, t_i) .$$

Learning an operator

$$\min_{F \in \mathcal{B}_0(\mathcal{Y}, \mathcal{X}), \Omega(F) < \infty} \{R_N(F) + \lambda \Omega(F)\} .$$

Theory

Is it a "good" algorithm in theory?

- To be investigated...
- See Srebro et al. (2004), Bach (2007) for preliminary results with the trace norm

Practice

Can we implement it? Does it work on real data?

- Optimization problem in the space of compact operators... but we show later that it boils down to a finite-dimensional optimization problem
- Promising results on real data

Theory

Is it a "good" algorithm in theory?

- To be investigated...
- See Srebro et al. (2004), Bach (2007) for preliminary results with the trace norm

Practice

Can we implement it? Does it work on real data?

- Optimization problem in the space of compact operators... but we show later that it boils down to a finite-dimensional optimization problem
- Promising results on real data

Theory

Is it a "good" algorithm in theory?

- To be investigated...
- See Srebro et al. (2004), Bach (2007) for preliminary results with the trace norm

Practice

Can we implement it? Does it work on real data?

- Optimization problem in the space of compact operators... but we show later that it boils down to a finite-dimensional optimization problem
- Promising results on real data

Theory

Is it a "good" algorithm in theory?

- To be investigated...
- See Srebro et al. (2004), Bach (2007) for preliminary results with the trace norm

Practice

Can we implement it? Does it work on real data?

- Optimization problem in the space of compact operators... but we show later that it boils down to a finite-dimensional optimization problem
- Promising results on real data

Theorem

If \hat{F} is a solution the problem:

$$\min_{F \in \mathcal{B}_2(\mathcal{Y}, \mathcal{X})} \left\{ R_N(F) + \lambda \sum_{i=1}^{\infty} \sigma_i(F)^2 \right\},$$

then it is necessarily in the linear span of $\{\mathbf{x}_i \otimes \mathbf{y}_i : i = 1, \dots, N\}$, i.e., it can be written as:

$$\hat{F} = \sum_{i=1}^N \alpha_i \mathbf{x}_i \otimes \mathbf{y}_i,$$

for some $\alpha \in \mathbb{R}^N$.

- $B_2(\mathcal{Y}, \mathcal{X})$ is isomorphic to the **RKHS** of the **tensor product kernel**:

$$k_{\otimes}((\mathbf{x}, \mathbf{y}), (\mathbf{x}', \mathbf{y}')) = \langle \mathbf{x}, \mathbf{x}' \rangle_{\mathcal{X}} \langle \mathbf{y}, \mathbf{y}' \rangle_{\mathcal{Y}},$$

by $f(\mathbf{x}, \mathbf{y}) = \langle \mathbf{x}, F\mathbf{y} \rangle_{\mathcal{X}}$. In particular,

$$\|f\|_{\mathcal{H}_{\otimes}}^2 = \|F\|^2 = \Omega(F).$$

- The problem is therefore a classical kernel method:

$$\min_{f \in \mathcal{H}_{\otimes}} \left\{ R_N(f) + \lambda \|f\|_{\otimes}^2 \right\},$$

so the classical representer theorem can be used. \square

A generalized representer theorem

Theorem

For **any spectral penalty function** $\Omega : \mathcal{B}_0(\mathcal{Y}, \mathcal{X}) \mapsto \mathbb{R}$, let the optimization problem:

$$\min_{F \in \mathcal{B}_0(\mathcal{Y}, \mathcal{X}), \Omega(F) < \infty} \{R_N(F) + \lambda \Omega(F)\} .$$

If the set of solutions is not empty, then there is a solution F in $\mathcal{X}_N \otimes \mathcal{Y}_N$, i.e., **there exists** $\alpha \in \mathbb{R}^{m_x \times m_y}$ **such that:**

$$F = \sum_{i=1}^{m_x} \sum_{j=1}^{m_y} \alpha_{ij} \mathbf{u}_i \otimes \mathbf{v}_j ,$$

where $(\mathbf{u}_1, \dots, \mathbf{u}_{m_x})$ and $(\mathbf{v}_1, \dots, \mathbf{v}_{m_y})$ form orthonormal bases of \mathcal{X}_N and \mathcal{Y}_N , respectively.

- For any operator $F \in \mathcal{B}_0(\mathcal{Y}, \mathcal{X})$, let

$$G = \Pi_{\mathcal{X}_N} F \Pi_{\mathcal{Y}_N},$$

where Π_U is the orthogonal projection onto U .

- Lemma: we can show that for all $i \geq 0$:

$$\sigma_i(G) \leq \sigma_i(F).$$

- Therefore $\Omega(G) \leq \Omega(F)$.
- On the other hand $R_N(G) = R_N(F)$.
- Consequently for any solution F we have another solution $G \in \mathcal{X}_N \otimes \mathcal{Y}_N$. \square

Theorem (cont.)

The coefficients α that define the solution by

$$F = \sum_{i=1}^{m_x} \sum_{j=1}^{m_y} \alpha_{ij} \mathbf{u}_i \otimes \mathbf{v}_j,$$

can be found by solving the following **finite-dimensional** optimization problem:

$$\min_{\alpha \in \mathbb{R}^{m_x \times m_y}, \Omega(\alpha) < \infty} R_N \left(\text{diag} \left(X \alpha Y^\top \right) \right) + \lambda \Omega(\alpha),$$

where $\Omega(\alpha)$ refers to the spectral penalty function applied to the matrix α seen as an operator from \mathbb{R}^{m_y} to \mathbb{R}^{m_x} , and X and Y denote any matrices that satisfy $K = XX^\top$ and $G = YY^\top$ for the two Gram matrices K and G of \mathcal{X}_N and \mathcal{Y}_N .

We obtain various algorithms by choosing:

- 1 A **loss function** (depends on the application)
- 2 A **spectral regularization** (that is amenable to optimization)
- 3 Two **Gram matrices** (aka kernel matrices)

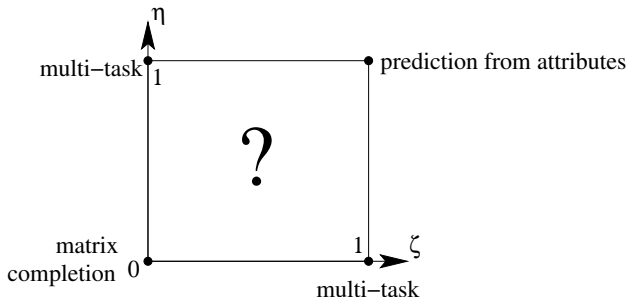
Both kernels and spectral regularization can be used to constrain the solution

A family of kernels

Taken $K_{\otimes} = K \times G$ with

$$\begin{cases} K = \eta K_{Attribute}^x + (1 - \eta) K_{Dirac}^x, \\ G = \zeta K_{Attribute}^y + (1 - \zeta) K_{Dirac}^y, \end{cases}$$

for $0 \leq \eta \leq 1$ and $0 \leq \zeta \leq 1$



Experiment

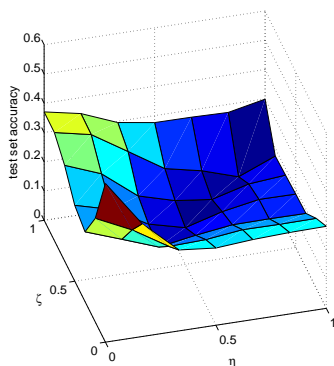
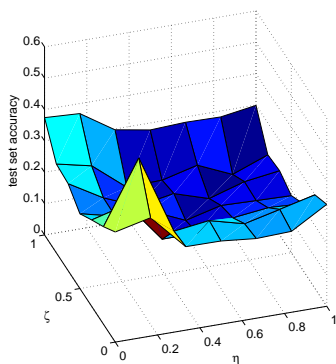
- Generate data $(\mathbf{x}, \mathbf{y}, z) \in \mathbb{R}^{f_X} \times \mathbb{R}^{f_Y} \times \mathbb{R}$ according to

$$z = \mathbf{x}^\top \mathbf{B} \mathbf{y} + \varepsilon$$

- Observe only $n_X < f_X$ and $n_Y < f_Y$ features
 - Low-rank assumption will find the missing features
 - Observed attributes will help the low-rank formulation to concentrate mostly on the unknown features
- Comparison of
 - Low-rank constraint without tracenorm (note that it requires regularization)
 - Trace-norm formulation (regularization is implicit)

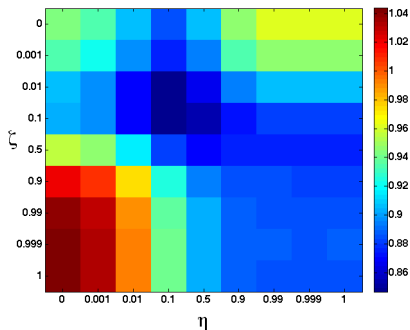
Simulated data: results

- Compare MSE
- Left: rank constraint (best: 0.1540), right: trace norm (best: 0.1522)



Movies

- MovieLens 100k database, ratings with attributes
- Experiments with 943 movies and 1,642 customers, 100,000 rankings in $\{1, \dots, 5\}$
- Train on a subset of the ratings, test on the rest
- error measured with MSE (best constant prediction: 1.26)



Conclusion

What we saw

- A general framework for CF with or without attributes
- A generalized representation theorem valid for any spectral penalty function
- A family of new methods;

Future work

- The bottleneck is often practical optimization. Online version possible.
- Automatic choice of the kernel

Reference

J. Abernethy, F. Bach, T. Evgeniou and J.-P. Vert, “A new approach to collaborative filtering: operator estimation with spectral constraint”, *technical report arXiv 0802-1430*, 2008.