**Case Study 4: Collaborative Filtering**

Collaborative Filtering
Matrix Completion
Alternating Least Squares

Machine Learning/Statistics for Big Data
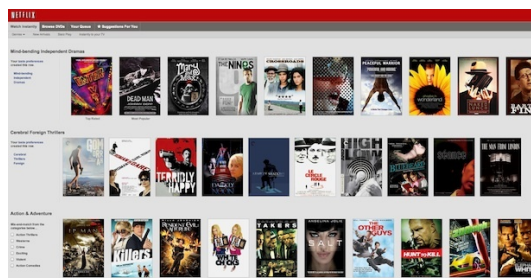CSE599C1/STAT592, University of Washington

Carlos Guestrin
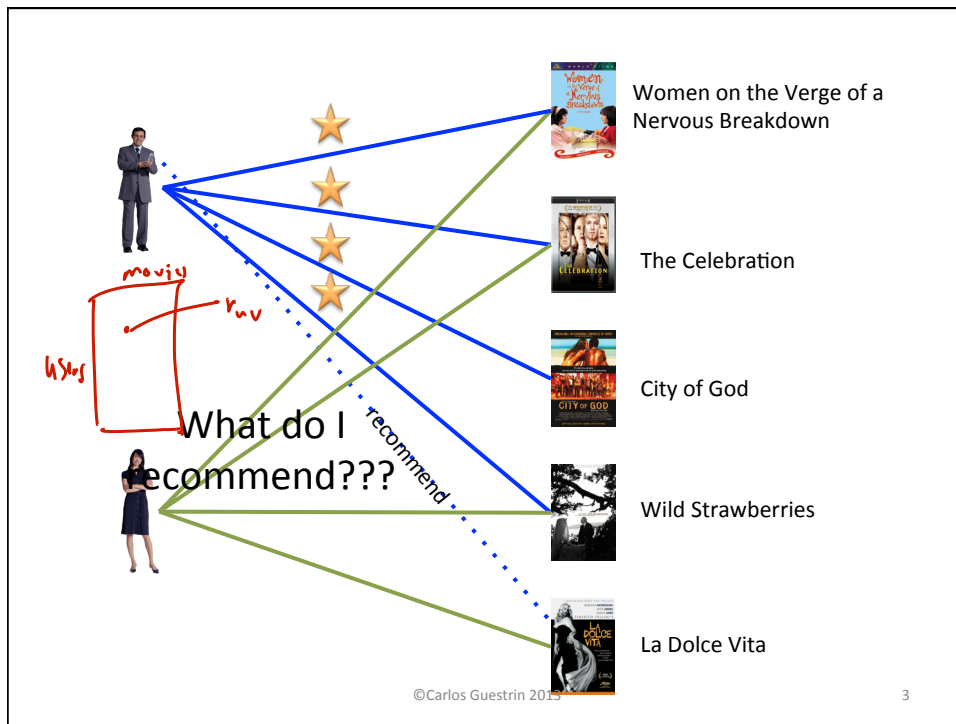
February 28th, 2013

1

---

# Collaborative Filtering

- **Goal:** Find movies of interest to a user based on movies watched by the user and others
- **Methods:** matrix factorization, GraphLab

2

Women on the Verge of a Nervous Breakdown

The Celebration

City of God

Wild Strawberries

La Dolce Vita

What do I recommend???

movie

$r_{uv}$

users

recommend

©Carlos Guestrin 2013

3

# Cold-Start Problem

- **Challenge:** Cold-start problem (new movie or user)
- **Methods:** use features of movie/user

$\phi(Skyfall) = \begin{pmatrix} 7 \\ 2 \\ 0 \\ \vdots \end{pmatrix}$ action romance

$\phi(FRWL) = \begin{pmatrix} 8 \\ 1 \\ 6 \\ \vdots \end{pmatrix}$

SKYFALL 007

IN THEATERS

©Carlos Guestrin 2013

4

2

# Netflix Prize

- Given 100 million ratings on a scale of 1 to 5, predict 3 million ratings to highest accuracy

Skyfall ⭐⭐⭐☆☆ 🚫 Not Interested

WÒNB ⭐☆☆☆☆ 🚫 Not Interested

FRWL ⭐⭐⭐⭐☆ 🚫 Not Interested

⭐⭐⭐⭐☆ 🚫 Not Interested  PoC

movies

users

8B params

- 17770 total movies
- 480189 total users
- Over 8 billion total ratings

- How to fill in the blanks?

& 100M obs

Figures from Ben Recht

©Carlos Guestrin 2013    5

---

# Matrix Completion Problem

movies

**X** =

users
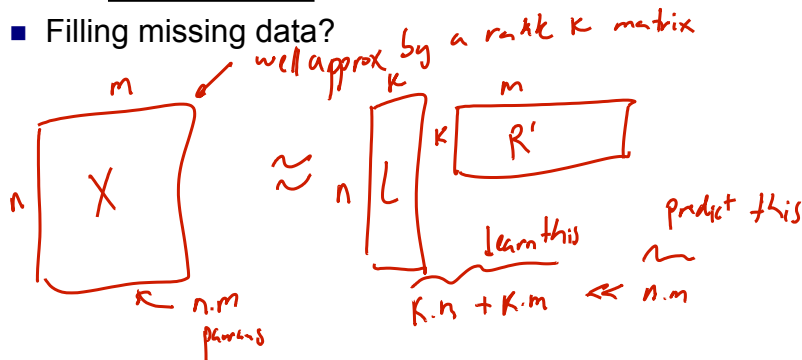
$X_{ij}$ known for black cells
$X_{ij}$ unknown for white cells
Rows index users
Columns index movies

- Filling missing data?

well approx by a rank K matrix

$X \approx L \; R'$

m

n   X

n.m params

predict this

learn this

$K \cdot n + K \cdot m \ll n \cdot m$

©Carlos Guestrin 2013    6

3

# Interpreting Low-Rank Matrix Completion (aka Matrix Factorization)



$r_{uv} \approx L_u \cdot R_v$

movie topic $i$ "action"

how much user $u$ likes topic $i$

how much movie $v$ is about topic $i$

7

---

# Matrix Completion via Rank Minimization

- Given observed values: $(u, v, r_{uv}) \in X$   Some $r_{uv} = ?$

- Find matrix $\hat{\Theta}$

- Such that: $\Theta_{uv} = r_{uv}$   $\forall r_{uv} \neq ?$

  fit $r_{uv} \neq ?$ perfectly

- But…

- Introduce bias:   Low rank

$$\min_{\Theta} \quad rank(\Theta)$$
$$\Theta_{uv} = r_{uv} \quad \forall r_{uv} \neq ?$$

- Two issues:   NP-hard   you can't hope to get exact matching

8

---

4

# Approximate Matrix Completion

- Minimize squared error:
  - (Other loss functions are possible)

$$\min_{\Theta} \sum_{r_{uv} \in X : r_{uv} \neq ?} (\Theta_{uv} - r_{uv})^2$$

- Choose rank *k:*

$$n \overset{m}{\Theta} = n \overset{k}{L} \cdot \overset{k}{{}_m R'}$$

- Optimization problem:

$$\min_{L,R} \sum_{r_{uv}} (L_u \cdot R_v - r_{uv})^2$$

non convex OPT problem ... local optima only

---

# Coordinate Descent for Matrix Factorization

$$\min_{L,R} \sum_{(u,v,r_{uv}) \in X : r_{uv} \neq ?} (L_u \cdot R_v - r_{uv})^2$$

- Fix movie factors, optimize for user factors

- First Observation:

$V_u \equiv$ set movies user u rated

$$\min_{L_1,..L_n} \sum_{(u,v,r_{uv})} (L_u \cdot R_v - r_{uv})^2 =$$

$$\min_{L_1,...L_n} \sum_u \sum_{v \in V_u} (L_u \cdot R_v - r_{uv})^2 = \quad \leftarrow \text{ indep opt problem per user}$$

$$\sum_u \min_{L_u} \sum_{v \in V_u} (L_u \cdot R_v - r_{uv})^2 \quad \leftarrow \text{ data parallel problem}$$
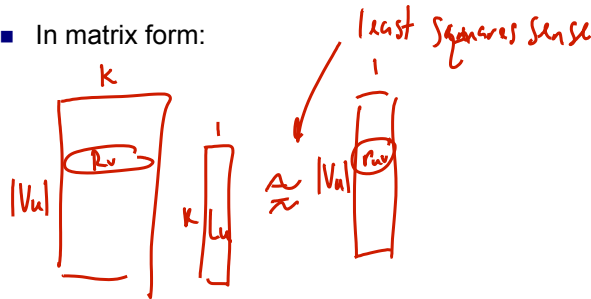
# Minimizing Over User Factors

- For each user $u$: $\displaystyle\min_{L_u}\sum_{v\in V_u}(L_u\cdot R_v - r_{uv})^2$

- In matrix form:

$least\ Squares\ Sense$



- Second observation: Solve by  - matrix inversion
  - gradient
  - ...

---

# Coordinate Descent for Matrix Factorization: Alternating Least-Squares

$$\min_{L,R}\sum_{(u,v,r_{uv})\in X: r_{uv}\neq ?}(L_u\cdot R_v - r_{uv})^2 + \lambda_u\|L\| + \lambda_v\|R\|$$

- Fix movie factors, optimize for user factors
  □ Independent least-squares over users   $\displaystyle\min_{L_u}\sum_{v\in V_u}(L_u\cdot R_v - r_{uv})^2 + \lambda_u\|L_u\|$

- Fix user factors, optimize for movie factors
  □ Independent least-squares over movies   $\displaystyle\min_{R_v}\sum_{u\in U_v}(L_u\cdot R_v - r_{uv})^2 + \lambda_v\|R_v\|$
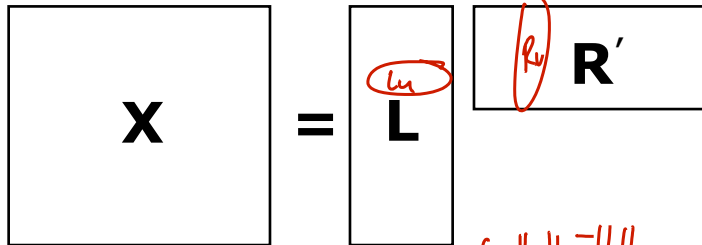
- System may be underdetermined: use regularization

- Converges to  local optima

# Effect of Regularization

$$\|A\|_F = \sqrt{\sum_{ij} A_{ij}^2}$$

$$\min_{L,R} \sum_{(u,v,r_{uv}) \in X : r_{uv} \neq ?} (L_u \cdot R_v - r_{uv})^2 + \lambda_u \|L\| + \lambda_v \|R\|$$

$$X = L \quad R'$$

if $\|\cdot\| \equiv \|\cdot\|_F^2$
each sub problem
uses $\|L_u\|_2^2 \rightarrow$ ridge regression

if $\|\cdot\| \equiv \|\cdot\|_1$
each subproblem $\|L_u\|_1$
$\Rightarrow$ solved by Lasso methods

13

---

# What you need to know…

- Matrix completion problem for collaborative filtering
- Over-determined -> low-rank approximation
- Rank minimization is NP-hard
- Minimize least-squares prediction for known values for given rank of matrix
  - Must use regularization
- Coordinate descent algorithm = "Alternating Least Squares"

14

# Case Study 4: Collaborative Filtering

## SGD for Matrix Completion
## Matrix-norm Minimization

Machine Learning/Statistics for Big Data
CSE599C1/STAT592, University of Washington

Carlos Guestrin

March 7th, 2013

15

---

# Stochastic Gradient Descent

$$\min_{L,R} F(L,R) = \min_{L,R} \frac{1}{2} \sum_{r_{uv}} (L_u \cdot R_v - r_{uv})^2 + \frac{\lambda_u}{2}||L||_F^2 + \frac{\lambda_v}{2}||R||_F^2$$

$\sum_u L_u \cdot L_u$

- Observe one rating at a time $r_{uv}^t$ $\quad \varepsilon_t = L_u^{(t)} \cdot R_v^{(t)} - r_{uv}$

- Gradient observing $r_{uv}$:

$$\frac{\partial F}{\partial L_u} = \varepsilon_t R_v + \lambda_u L_u$$
$$\frac{\partial F}{\partial R_v} = \varepsilon_t L_u + \lambda_v R_v$$

$$\nabla F_t = \begin{bmatrix} \varepsilon_t R_v + \lambda_u L_u \\ \varepsilon_t L_u + \lambda_v R_v \end{bmatrix}$$

- Updates: step size $\eta_t$. $\begin{bmatrix} L \\ R \end{bmatrix} \leftarrow \begin{bmatrix} L \\ R \end{bmatrix} - \eta_t \nabla F_t$

fast & easy to implement

$$\begin{bmatrix} L_u^{(t+1)} \\ R_v^{(t+1)} \end{bmatrix} \leftarrow \begin{bmatrix} (1 - \eta_t \lambda_u) L_u^{(t)} - \eta_t \varepsilon_t R_v^{(t)} \\ (1 - \eta_t \lambda_v) R_v^{(t)} - \eta_t \varepsilon_t L_u^{(t)} \end{bmatrix}$$

16

8

# Local Optima v. Global Optima

- We are solving:

$$\min_{L,R} \sum_{r_{uv}} (L_u \cdot R_v - r_{uv})^2 + \lambda_u ||L||_F^2 + \lambda_v ||R||_F^2$$

- We (kind of) wanted to solve:

$$\min_{\theta} \; rank(\theta)$$
$$\theta_{uv} = r_{uv} \qquad \forall r_{uv} \in X, \; r_{uv} \pm ?$$

- Which is NP-hard…
  - How do these things relate???

---

# Eigenvalue Decompositions for PSD Matrices

- Given a (square) symmetric positive semidefinite matrix: $\theta \succeq 0$
  - Eigenvalues: $\lambda_1, \cdots, \lambda_d \geq 0$ $\qquad \lambda = (\lambda_1, \cdots, \lambda_d)$
- Thus rank is:

$$|\{\lambda_i : \lambda_i > 0\}| \equiv rank(\theta) \equiv ||\lambda||_0$$

- Approximation:

$$||\lambda||_0 \approx ||\lambda||_1 = \sum_{i=1}^{d} |\lambda_i| \overset{PSD}{=} \sum_{i=1}^{d} \lambda_i \checkmark \quad L_1 \text{ norm is} \quad sum \; of \; eigen \; values$$

- Property of trace:

$$trace(\theta) = \sum_{i=1}^{d} \lambda_i$$

- Thus, approximate rank minimization by:

$$\min_{\theta} \; trace(\theta)$$
$$\theta_{uv} = r_{uv}$$
$$\theta \succeq 0$$

# Generalizing the Trace Trick

- Non-square matrices ain't got no trace

- For (square) positive definite matrices, matrix factorization: $\lambda_i \geq 0$

$$d\hat{\Theta} = P \Lambda P^{-1} \qquad diag(\lambda)$$

- For rectangular matrices, singular value decomposition: $e.g.\ n \geq m$

$$n\ \boxed{\Theta}\ = \ n\ \boxed{U}\ \ n\ \boxed{\Sigma}\ \ m\ \boxed{V'}$$

diagonal matrix
entries
$\sigma_i(\theta) \geq 0$  ith singular value

- Nuclear norm:

$$\| \Theta \|_* = \sum_{i=1}^{m} \sigma_i(\theta) \qquad \begin{array}{l} \min_{\theta} \|\theta\|_* \\ \Theta_{uv} = r_{uv} \end{array} \quad \begin{array}{l} \text{convex} \\ \text{problem} \end{array}$$

---

# Nuclear Norm Minimization

- Optimization problem:

$$\min_{\theta} \|\theta\|_* $$
$$\Theta_{uv} = r_{uv}$$

- Possible to relax equality constraints:

$$\min_{\theta} \sum_{ruv} (\Theta_{uv} - r_{uv})^2 + \lambda \|\theta\|_*$$

- Both are convex problems!
  (solved by <u>semidefinite programming</u>)

# Analysis of Nuclear Norm

- Nuclear norm minimization is a convex relaxation of rank minimization problem:

*NP-hard*

*convex relaxation with a poly time solution*

$$\min_{\Theta} rank(\Theta) \qquad\qquad \min_{\Theta} ||\Theta||_*$$

$$r_{uv} = \Theta_{uv}, \forall r_{uv} \in X, r_{uv} \neq ? \qquad\qquad r_{uv} = \Theta_{uv}, \forall r_{uv} \in X, r_{uv} \neq ?$$

- Theorem [Candes, Recht '08]:
  - □ If there is a true matrix of rank *k*,
  - □ And, we observe at least

*original problem n·m entries*

*rank*

*need about $k n^{1.2}$*

*constant*

$$C\ k\ n^{1.2}\log n$$

*we have $Kn + km$ params*

*$n \geq m$*

   random entries of true matrix

  - □ Then true matrix is recovered exactly with high probability with convex nuclear norm minimization!
    - Under certain conditions

©Carlos Guestrin 2013

21

---

# Nuclear Norm Minimization versus Direct (Bilinear) Low Rank Solutions

- Nuclear norm minimization:

*$\Theta^* = args$*

$$\min_{\Theta} \sum_{r_{uv}} (\Theta_{uv} - r_{uv})^2 + \lambda ||\Theta||_* \qquad (*)$$

*Convex, global OPT = close to truth*

  - □ Annoying because: — $\Theta$ *very large* (8B entries in Netflix)
    - *SDP solvers are very slow (but poly time)*

- Instead: $\min_{L,R} \sum_{r_{uv}} (L_u \cdot R_v - r_{uv})^2 + \underline{\lambda_u ||L||_F^2 + \lambda_v ||R||_F^2} \qquad (**)$

*≈ 10-100M params, solvers very fast*

  - □ Annoying because: — *many local optima*

  - □ But $||\Theta||_* = \inf \left\{ \min_{L,R} \frac{1}{2}||L||_F^2 + \frac{1}{2}||R||_F^2 : \Theta = LR' \right\}$

    - So *second prob. nonconvex approx. to first*
    - And *if pick rank of L & R to be slightly higher rank($\Theta^*$), local optima of (**) are global optima of (*)*
    - Under certain conditions [Burer, Monteiro '04]

©Carlos Guestrin 2013

22

---

11

# What you need to know…

- Stochastic gradient descent for matrix factorization

- Norm minimization as convex relaxation of rank minimization
  - Trace norm for PSD matrices
  - Nuclear norm in general

- Intuitive relationship between nuclear norm minimization and direct (bilinear) minimization

**23**

---

# Case Study 4: Collaborative Filtering

## Nonnegative Matrix Factorization
## Projected Gradient

Machine Learning/Statistics for Big Data
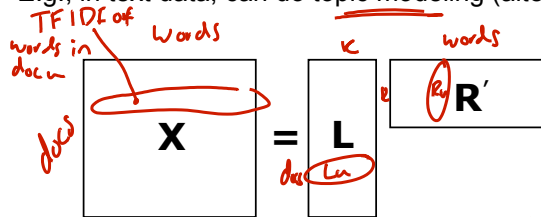CSE599C1/STAT592, University of Washington

Carlos Guestrin

March 7$^{th}$, 2013

**24**

# Matrix factorization solutions can be unintuitive…

- Many, many, many applications of matrix factorization

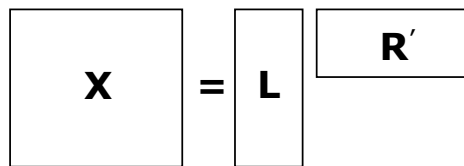- E.g., in text data, can do topic modeling (alternative to LDA):

TF IDF of words in doc $u$

words $k$ words

$R'$

$X$ $=$ $L$

docs $R_v$

doc $L_u$

- Would like:

$L_u$: how much a doc is about each topic

$R_v$: how much a word contributes to each topic

- But… Standard matrix factorization: $L_u, R_v$ can be negative

---

# Nonnegative Matrix Factorization

$R'$

$X$ $=$ $L$

- Just like before, but

$$\min_{L \geq 0, R \geq 0} \sum_{r_{uv}} (L_u \cdot R_v - r_{uv})^2 + \lambda_u ||L||_F^2 + \lambda_v ||R||_F^2$$
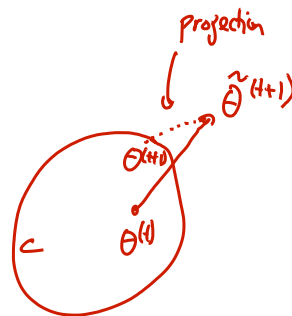
nonnegative $L, R$

- Constrained optimization problem
  - Many, many, many, many solution methods… we'll check out a simple one

# Projected Gradient

- Standard optimization:
  - Want to minimize: $\min_{\Theta} f(\Theta)$
  - Use gradient updates:
    $$\Theta^{(t+1)} \leftarrow \Theta^{(t)} - \eta_t \nabla f(\Theta^{(t)})$$
- Constrained optimization:
  - Given convex set *C* of feasible solutions
  - Want to find minima within *C*: $\min_{\Theta} f(\Theta)$
    $$\Theta \in \mathcal{C}$$
- Projected gradient:
  - Take a gradient step (ignoring constraints):
    $$\tilde{\Theta}^{(t+1)} \leftarrow \Theta^{(t)} - \eta_t \nabla f(\Theta^{(t)})$$
  - Projection into feasible set:
    $$\Pi_c(\theta) \equiv \operatorname*{argmin}_{\beta \in C} \|\theta - \beta\|_2^2 \Big\} \text{ often easy to compute (always convex)}$$
    $$\Theta^{(t+1)} = \Pi_c\left(\tilde{\Theta}^{(t+1)}\right)$$

*Projection*

---

# Projected Stochastic Gradient Descent for Nonnegative Matrix Factorization

$$\min_{L \geq 0, R \geq 0} \frac{1}{2} \sum_{r_{uv}} (L_u \cdot R_v - r_{uv})^2 + \frac{\lambda_u}{2}\|L\|_F^2 + \frac{\lambda_v}{2}\|R\|_F^2$$

- Gradient step observing $r_{uv}$ ignoring constraints:
$$\begin{bmatrix} \tilde{L}_u^{(t+1)} \\ \tilde{R}_v^{(t+1)} \end{bmatrix} \leftarrow \begin{bmatrix} (1 - \eta_t \lambda_u)L_u^{(t)} - \eta_t \epsilon_t R_v^{(t)} \\ (1 - \eta_t \lambda_v)R_v^{(t)} - \eta_t \epsilon_t L_u^{(t)} \end{bmatrix}$$

- Convex set: $L_u \geqslant 0 \quad ; \quad R_v \geqslant 0 \quad \forall u,v$
- Projection step:

$$\Pi_c(\theta) = \operatorname*{argmin}_{\beta \in C} \|\theta - \beta\|_2^2 \leftarrow \text{totally indep prob per dimension}$$

Single dim
$$= \operatorname*{argmin}_{\beta \geqslant 0} (\theta - \beta)^2$$
$$= \begin{cases} \theta, & \text{if } \theta \geqslant 0 \\ 0, & \text{if } \theta < 0 \end{cases} = (\theta)_+$$

$$\begin{bmatrix} L^{(t+1)} \\ R^{(t+1)} \end{bmatrix} \leftarrow \begin{bmatrix} \tilde{L}^{(t+1)} \\ \tilde{R}^{(t+1)} \end{bmatrix}$$

set all neg coords to zero, universe on our side!!

# What you need to know…

- In many applications, want factors to be nonnegative

- Corresponds to constrained optimization problem

- Many possible approaches to solve, e.g., projected gradient

29