It's a Small(-Rank) World, After All

Geoff Gordon SELECT Lab Machine Learning Department Joint work with Ajit Singh, Francisco Pereira, Nick Roy, Byron Boots, Sajid Siddiqi

Principal Components Analysis



PCA: the equation





rows of V span the lowrank space

Data matrix X

Compressed matrix U

Examples: topic analysis

CHAMPION PRODUCTS <CH> APPROVES STOCK SPLIT

The watching and the second railers in the fame are the process of the process of the process of the second s

ROCHESTER, N.Y., Feb 26 - Champion Products Inc said its board of directors approved a two-for-one stock split of its common shares for shareholders of record as of April 1, 1987.

The company also said its board voted to recommend to shareholders at the annual meeting April 23 an increase in the authorized capital stock from five mln to 25 mln shares. Reuter

Data matrix: "bag of words"

mit interteditor

words



Result of factoring







Basis vectors often called "topics" Typical topic: corn, future, price, weather, market, ...

Example: images



face images from Groundhog Day, extracted by Cambridge face DB project

Data matrix

LTOT WAST'S mon - Sector Destadiates in 'De Low day the transmission

pixels

STREET & . ANYES



Result of factoring



Eigenfaces



image credit: AT&T Labs Cambridge

Example: collaborative filtering

movies



And so forth...



 bibliometrics, spectral clustering, speech recognition, PageRank, ...

image credit: <u>http://www.cheswick.com/ches/</u>

If you have a hammer



This talk: some more nails

and a state of the second and the se



image credit: Wikipedia

...and some bigger hammers



image credit: Wikipedia

Interesting nail #1: movies



Can we learn to predict future frames?

Data matrix: pixel sequences



pixels from the future

Result of factoring





Basis weights ut encapsulate state of system at step t After factoring, we can try to learn to predict ut+1 from ut

Movie example

STREET

the surperter the states in the Low did the the second with the



Movie example

N 100 6 1

and in the second of the second



Interesting nail #2: histograms



Interesting nail #2: histograms



Compressing a set of histograms



Why compress histograms?

States of the



Robot beliefs

Why compress histograms?

 Single histogram = test performance distribution (of multiple students on one test, or one student on multiple tests)

Interesting nail #3: relations

genres

 X_I

 \mathfrak{X}_2

X3

movies

users



does movie i have genre j?

 \mathcal{X}_{n}

Bigger hammers

STREET & ANTER

Nonlinearity
Relational factorization

Bigger hammer: nonlinearity



And its brother: non-normality



image credit: risklab.ch

Where do we need nonlinearity?

STREET & . BATES

• Bounds:

Images: can't have pixel >1 or <0
Word counts: can't be negative

Where do we need nonlinearity?

Consistency of interpolation
 Suppose x₁, x₂ are samples
 Linearity means x₁ + λ(x₂ - x₁) should make sense for any λ

Consistency of interpolation



training data predicted image

Where do we need non-normality?

Uneven variance
Error of ±3 in count of 5 vs. 50 vs. 500
Asymmetric errors
E.g., counts often have long tail

Uneven variance, asymmetric

errors



Learned (PCA)

Generalizing PCA

strend & - Advera

- In PCA, had $X_{ij} \approx U_i \cdot V_j$ • What if • $X_{ij} \approx exp(U_i \cdot V_j)$ • $X_{ij} \approx logit(U_i \cdot V_j)$
 - 0 ...

Generalizing PCA

STREET & . ANTERI

- In PCA, had X_{ij} ~ Normal(μ), μ = U_i · V_j
 What if
 - $X_{ij} \sim Poisson(\mu)$
 - $X_{ij} \sim Binomial(p)$
 - o ...?

Exponential family review

Exponential family of distributions: P(X | θ) = P₀(X) exp(Xθ – g(θ))
g(θ) is always strictly convex, differentiable on interior of domain
means g' is strictly monotone (in 1D)

Exponential family review

Exponential family PDF: P(X | θ) = P₀(X) exp(Xθ – g(θ))
g'(θ) = E(X | θ)
g' & (g')⁻¹ = "link function"
E(X | θ) = "expectation parameter"

STREET & . ANTES

Normal(mean)
Poisson(log rate)
Binomial(log odds)

Solving both problems at once

Let P(X | θ) be an exponential family with natural parameter θ
Predict X_{ij} ~ P(X | θ_{ij}), where θ_{ij} = U_i · V_j
e.g., in Poisson, E(X_{ij}) = exp(θ_{ij})
e.g., in Binomial, E(X_{ij}) = logit(θ_{ij})

More precisely,

 $\max_{U,V} \sum_{ij} \log P(X_{ij} | \theta_{ij}) + \log P(U) + \log P(V)$ s.t. $\theta_{ij} = U_i \cdot V_j$

 "Generalized linear" or "exponential family" PCA

• all P(...) terms are exponential families

• analogy to GLMs

[Collins et al, 2001] [Gordon, 2002] [Roy & Gordon, 2005]

Theorem

In GL PCA, finding U which maximizes likelihood (holding V fixed) is a convex optimization problem

And, finding best V (holding U fixed) is a convex problem

Further, Hessian is block diagonal

• So, an efficient and effective optimization algorithm: alternately improve U and V

Proof

WHEER A. ANYER

• Exponential family PDF: $P(X \mid \theta) = P_0(X) \exp(X\theta - g(\theta))$ • $\sum \log P(X_{ij} \mid \theta_{ij}) =$ $\sum [\log P_0(X_{ij}) - X_{ij}\theta_{ij} - g(\theta_{ij})] =$ $\sum [\log P_0(X_{ij}) - X_{ij}(U_i^TV_j) - g(U_i^TV_j)]$

Special cases

- PCA, probabilistic PCA
- k-means clustering
- Independent components analysis (ICA)
- Poisson PCA
- Max-margin matrix factorization (MMMF)
- Almost: pLSI, pHITS, NMF





infort Thereofer

NY 184 8 1











Robot beliefs example



Original







Learned (Poisson PCA)