The Log-Determinant Divergence and its Applications

Inderjit S. Dhillon University of Texas at Austin

Householder Symposium XVII Zeuthen, Germany June 3, 2008

Joint work with Jason Davis, Kristen Grauman, Prateek Jain, Brian Kulis, Suvrit Sra and Joel Tropp

Motivating Application: Image Search

- Thousands to millions of pixels in an image
- 3,000-30,000 human recognizable object categories
- Billions of images indexed by Google Image Search
- How can images be represented?
 - Global representations: One vector per image pixel intensities, color histograms, etc.
 - Local representations: Detect distinctive interest points and extract descriptors invariant to scale, translation, rotation, illumination, etc.

Image Search: Pyramid Match Kernel

- Use Local Representations to represent each image as a set of fixed-dimensional vectors
- Used the pyramid match kernel of [Grauman and Darrell, 2007] to compute approximate matching between two images
 - Place multidimensional, multi-resolution grid over point sets
 - Compute intersection of multi-dimensional histograms at multiple resolutions
 - Compute match between image *u* and *v* as:

$$\mathcal{K}(u,v)=\sum_{i=0}^L w_i N_i(u,v), \quad ext{where} \quad w_i>w_{i+1}>0,$$

伺 ト イ ヨ ト イ ヨ ト

and N_i is number of newly matched pairs at level i

• Can show that K is a positive definite matrix (kernel function)

Example query results

Query

Top three retrieved images



伺 ト イヨト イヨト

Example (bad) query results

Query

Top three retrieved images



□ ▶ < □ ▶ < □</p>

Inderjit S. Dhillon University of Texas at Austin The Log-Determinant Divergence and its Applications

- 4 同 6 4 日 6 4 日 6

э

Example: Image Search Is Not Perfect...

charlie van Ioan - Google Bilder

Web Bilder Maps News Shopping Mail Mehr V

Anmelden

Erweiterte Bildsuche Einstellungen



Bilder



Charles Van Loan 450 x 600 - 241k - aif www.cs.comell.edu



Mehr von www.cs. cornell.edu]

Top to bottom: Joe Van Loan. ... 256 x 256 - 50k - gif 438 x 594 - 19k - jpg www.cs.comell.edu inkspots.ca



... guitar; Joe Van Loan. 1st tenor: ... 247 x 249 - 7k - jpg inkspots.ca



by Charles E. Van Loan 400 x 400 - 43k - ipg gaslight.mtroyal.ca

Ergebnisse 1 - 20 von ungefähr 163.000 für charlie van Ioan. (0,05 Sekunden) Anzeigen: Alle Größen - Extra groß - groß - mittel - klein





... Charles Van Loan, Dennis

West. ...

308 x 250 - 29k - jpg

www.northerninitiatives.com

405 x 500 - 19k - jpg

Patricia Van Loan, Manotick ? -2007

The Log-Determinant Divergence and its Applications



Van Dooren Charlie Van Loan. ... 489 x 312 - 27k - jpg www.cs.umd.edu



... Charles F. Van Loan 501 x 798 - 62k - ing

Inderjit S. Dhillon University of Texas at Austin



Van Loan and Acord 216 x 388 - 10k - ipa gaslight.mtroyal.ca

Matrix Computation THERP EDUTION

Gene H. Golub and Charles F.

Van

Improving the Kernel

• The Matrix Nearness Problem:

min loss(K, K_0) $(\mathbf{e}_i - \mathbf{e}_j)^T K(\mathbf{e}_i - \mathbf{e}_j) \le u$ if $(i, j) \in S$ [similarity constraints] $(\mathbf{e}_i - \mathbf{e}_j)^T K(\mathbf{e}_i - \mathbf{e}_j) \ge \ell$ if $(i, j) \in D$ [dissimilarity constraints] $K \succeq 0$

- Learn kernel matrix K that is "close" to the baseline kernel matrix K_0
- Other linear constraints on K are possible
- Constraints can arise from various scenarios
 - Unsupervised: Click-through feedback
 - Semi-supervised: must-link and cannot-link constraints
 - Supervised: points in the same class have "small" distance, etc.

通 と イ ヨ と イ ヨ と

Improving the Kernel

• The Matrix Nearness Problem:

min loss(K, K_0) $(\mathbf{e}_i - \mathbf{e}_j)^T K(\mathbf{e}_i - \mathbf{e}_j) \le u$ if $(i, j) \in S$ [similarity constraints] $(\mathbf{e}_i - \mathbf{e}_j)^T K(\mathbf{e}_i - \mathbf{e}_j) \ge \ell$ if $(i, j) \in D$ [dissimilarity constraints] $K \succeq 0$

- Learn kernel matrix K that is "close" to the baseline kernel matrix K_0
- Other linear constraints on K are possible
- Constraints can arise from various scenarios
 - Unsupervised: Click-through feedback
 - Semi-supervised: must-link and cannot-link constraints
 - Supervised: points in the same class have "small" distance, etc.
- QUESTION: What should "loss" be?

ヨッ イヨッ イヨッ

Improving the Kernel

• The Matrix Nearness Problem:

min loss(K, K_0) $(\mathbf{e}_i - \mathbf{e}_j)^T K(\mathbf{e}_i - \mathbf{e}_j) \le u$ if $(i, j) \in S$ [similarity constraints] $(\mathbf{e}_i - \mathbf{e}_j)^T K(\mathbf{e}_i - \mathbf{e}_j) \ge \ell$ if $(i, j) \in D$ [dissimilarity constraints] $K \succeq 0$

- Learn kernel matrix K that is "close" to the baseline kernel matrix K_0
- Other linear constraints on K are possible
- Constraints can arise from various scenarios
 - Unsupervised: Click-through feedback
 - Semi-supervised: must-link and cannot-link constraints
 - Supervised: points in the same class have "small" distance, etc.
- QUESTION: What should "loss" be?
- We use "loss" to be the LogDet Divergence

通 と イ ヨ と イ ヨ と

What is the LogDet Divergence?

同 ト イ ヨ ト イ ヨ ト

• Frobenius Distance:

$$D_{Frob} = \|X - Y\|_F$$

・ 同 ト ・ ヨ ト ・ ヨ ト

э

• Frobenius Distance:

$$D_{Frob} = \|X - Y\|_F$$

• LogDet Divergence:

$$D_{\ell d}(X,Y) = \operatorname{trace}(XY^{-1}) - \log \det(XY^{-1}) - d,$$

3 N

• Frobenius Distance:

$$D_{Frob} = \|X - Y\|_F$$

• LogDet Divergence:

$$D_{\ell d}(X, Y) = \operatorname{trace}(XY^{-1}) - \log \det(XY^{-1}) - d,$$

= $\operatorname{trace}(Y^{-1/2}XY^{-1/2}) - \log \det(Y^{-1/2}XY^{-1/2}) - d,$

3 N

• Frobenius Distance:

$$D_{Frob} = \|X - Y\|_F$$

• LogDet Divergence:

$$D_{\ell d}(X, Y) = \operatorname{trace}(XY^{-1}) - \log \det(XY^{-1}) - d,$$

= $\operatorname{trace}(Y^{-1/2}XY^{-1/2}) - \log \det(Y^{-1/2}XY^{-1/2}) - d,$
= $\operatorname{trace}(Y^{-1/2}XY^{-1/2} - \log Y^{-1/2}XY^{-1/2} - I),$

3 N

• Frobenius Distance:

$$D_{Frob} = \|X - Y\|_F$$

LogDet Divergence:

$$D_{\ell d}(X, Y) = \operatorname{trace}(XY^{-1}) - \log \det(XY^{-1}) - d,$$

= $\operatorname{trace}(Y^{-1/2}XY^{-1/2}) - \log \det(Y^{-1/2}XY^{-1/2}) - d,$
= $\operatorname{trace}(Y^{-1/2}XY^{-1/2} - \log Y^{-1/2}XY^{-1/2} - I),$
= $\sum_{i} \sum_{j} (\mathbf{v}_{i}^{T}\mathbf{u}_{j})^{2} \left(\frac{\lambda_{i}}{\theta_{j}} - \log \frac{\lambda_{i}}{\theta_{j}} - 1\right),$

where $X = V \Lambda V^T$ and $Y = U \Theta U^T$

$$D_{\ell d}(X,Y) \;\;=\;\; {
m trace}(Y^{-1/2}XY^{-1/2}) - \log \det(Y^{-1/2}XY^{-1/2}) - d,$$

- Can be used as a measure of distance
 - Positive, and zero iff X = Y
 - But not symmetric, and triangle inequality does not hold
- Convex in first argument (not in second)
- Pythagorean Property holds: Given Y and a convex set Ω ,

 $D_{\ell d}(X,Y) \geq D_{\ell d}(X,P_{\Omega}(Y)) + D_{\ell d}(P_{\Omega}(Y),Y), \quad \text{holds for all } X \in \Omega$

• Definition can be extended to semi-definite matrices

LogDet Divergence: Scale Invariance

$$D_{\ell d}(X,Y) = \operatorname{trace}(Y^{-1/2}XY^{-1/2}) - \log \det(Y^{-1/2}XY^{-1/2}) - d,$$

Scale-invariance

$$D_{\ell d}(X, Y) = D_{\ell d}(\alpha X, \alpha Y), \quad \alpha \ge 0$$

• In fact, for any invertible M

$$D_{\ell d}(X,Y) = D_{\ell d}(M^T X M, M^T Y M)$$

• In particular,

$$D_{\ell d}(X,Y) = D_{\ell d}(Y^{-1/2}XY^{-1/2},I)$$

$$egin{array}{rcl} D_{\ell d}(X,I) &=& ext{trace}(X) - \log \det(X) - d, \ &=& \displaystyle{\sum_{i=1}^d \left(\lambda_i - \log \lambda_i - 1
ight)} \end{array}$$

<ロ> <同> <同> < 同> < 同>

э

$$egin{array}{rcl} D_{\ell d}(X,I) &=& ext{trace}(X) - \log \det(X) - d, \ &=& \displaystyle\sum_{i=1}^d \left(\lambda_i - \log \lambda_i - 1
ight) \end{array}$$

• Now, $x - \log x \ge 1$ with equality at x = 1

▲圖 → ▲ 国 → ▲ 国 → …

3

$$\begin{array}{lll} D_{\ell d}(X,I) &=& \operatorname{trace}(X) - \log \det(X) - d, \\ &=& \displaystyle \sum_{i=1}^d \left(\lambda_i - \log \lambda_i - 1\right) \end{array}$$

• Now, $x - \log x \ge 1$ with equality at x = 1

• Also, $x - \log x \ge \log x + 1 - \log 4$ with equality at x = 2

伺 と く き と く き と

-

$$egin{array}{rcl} D_{\ell d}(X,I) &=& ext{trace}(X) - \log \det(X) - d, \ &=& \displaystyle{\sum_{i=1}^d \left(\lambda_i - \log \lambda_i - 1
ight)} \end{array}$$

• Now, $x - \log x \ge 1$ with equality at x = 1

- Also, $x \log x \ge \log x + 1 \log 4$ with equality at x = 2
- Letting, $\lambda_1 \geq \lambda_2 \cdots \geq \lambda_d > 0$

$$\begin{array}{rcl} D_{\ell d}(X,I) & \geq & (\log \lambda_1 + 1 - \log 4) - (\log \lambda_d + 1), \\ \Longrightarrow & \operatorname{cond}(X) & \leq & 4 \exp D_{\ell d}(X,I) \end{array}$$

Thus, LogDet yields an upper bound on the condition number

伺 ト イ ヨ ト イ ヨ ト

Where does LogDet occur?

Inderjit S. Dhillon University of Texas at Austin The Log-Determinant Divergence and its Applications

同 ト イ ヨ ト イ ヨ ト

-

The Gaussian Connection — Maximum Likelihood

• Suppose $\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_m$ are drawn from:

$$p(\mathbf{x}|\boldsymbol{\mu},\boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{d/2} (\det \boldsymbol{\Sigma})^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})\right\}$$

伺 と く き と く き と

3

The Gaussian Connection — Maximum Likelihood

• Suppose $\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_m$ are drawn from:

$$p(\mathbf{x}|\boldsymbol{\mu},\boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{d/2} (\mathsf{det}\boldsymbol{\Sigma})^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^{\mathsf{T}}\boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})\right\}$$

Log-Likelihood:

$$L(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \prod_{i=1}^{m} p(\mathbf{x}_i | \boldsymbol{\mu}, \boldsymbol{\Sigma})$$

$$\propto \exp \left\{ -\frac{m}{2} \left(D_{\ell d}(\boldsymbol{\Sigma}^{-1}, \bar{\mathbf{S}}^{-1}) + (\bar{\boldsymbol{\mu}} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\bar{\boldsymbol{\mu}} - \boldsymbol{\mu}) \right) \right\},$$

where $\bar{\boldsymbol{\mu}} = \frac{1}{m} \sum_{i=1}^{m} \mathbf{x}_i$ and $\bar{\mathbf{S}} = \frac{1}{m} \sum_{i=1}^{m} (\mathbf{x}_i - \bar{\boldsymbol{\mu}}) (\mathbf{x}_i - \bar{\boldsymbol{\mu}})^T$

同 ト イヨ ト イヨ ト ヨ うくや

The Gaussian Connection — Maximum Likelihood

• Suppose $\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_m$ are drawn from:

$$p(\mathbf{x}|oldsymbol{\mu},oldsymbol{\Sigma}) = rac{1}{(2\pi)^{d/2}(\mathsf{det}oldsymbol{\Sigma})^{1/2}} \exp\left\{-rac{1}{2}(\mathbf{x}-oldsymbol{\mu})^T oldsymbol{\Sigma}^{-1}(\mathbf{x}-oldsymbol{\mu})
ight\}$$

Log-Likelihood:

$$\begin{split} \mathcal{L}(\boldsymbol{\mu},\boldsymbol{\Sigma}) &= \prod_{i=1}^{m} p(\mathbf{x}_{i} | \boldsymbol{\mu},\boldsymbol{\Sigma}) \\ &\propto & \exp \Big\{ -\frac{m}{2} \left(D_{\ell d}(\boldsymbol{\Sigma}^{-1},\bar{\mathbf{S}}^{-1}) + (\bar{\boldsymbol{\mu}}-\boldsymbol{\mu})^{T} \boldsymbol{\Sigma}^{-1} (\bar{\boldsymbol{\mu}}-\boldsymbol{\mu}) \right) \Big\}, \end{split}$$

where $\bar{\boldsymbol{\mu}} = \frac{1}{m} \sum_{i=1}^{m} \mathbf{x}_i$ and $\bar{\mathbf{S}} = \frac{1}{m} \sum_{i=1}^{m} (\mathbf{x}_i - \bar{\boldsymbol{\mu}}) (\mathbf{x}_i - \bar{\boldsymbol{\mu}})^T$

• Thus, $\bar{\mu}$ and $\bar{\mathbf{S}}$ are the maximum likelihood estimates of the mean & covariance matrix, respectively

Further Connections in Statistics

• Wishart Distribution — Given *m* samples from a Gaussian distribution, the pdf of the sample covariance matrix may be written as:

$$p(\mathbf{S},m) \propto \exp\left\{-rac{m}{2}D_{\ell d}(\mathbf{\Sigma^{-1}},\mathbf{S}^{-1}) - rac{d+1}{2}\log\det\mathbf{S}
ight\}$$

• Differential Relative Entropy between two Multivariate Gaussians:

$$\int p(\mathbf{x}|\boldsymbol{\mu},\boldsymbol{\Sigma}_0) \log\left(\frac{p(\mathbf{x}|\boldsymbol{\mu},\boldsymbol{\Sigma}_0)}{p(\mathbf{x}|\boldsymbol{\mu},\boldsymbol{\Sigma})}\right) d\mathbf{x} = \frac{1}{2} D_{\ell d}(\boldsymbol{\Sigma},\boldsymbol{\Sigma}_0)$$

• [James and Stein, 1961] LogDet divergence is known as Stein's loss in the statistics community

Quasi-Newton Optimization

- LogDet Divergence arises in the BFGS and DFP updates
 - Quasi-Newton methods
 - Approximate Hessian of the function to be minimized
- [Fletcher, 1991] BFGS update can be shown to optimize:

$$\begin{array}{ll} \min_{B} & D_{\ell d}(B,B_{t})\\ \text{subject to} & B \ s_{t}=y_{t} \ (\text{``Secant Equation''}) \end{array}$$

•
$$s_t = x_{t+1} - x_t$$
, $y_t = \nabla f_{t+1} - \nabla f_t$

• Closed-form solution:

$$B_{t+1} = B_t - \frac{B_t s_s s_t^T B_t}{s_t^T B_t s_t} + \frac{y_t y_t^T}{s_t^T y_t}$$

• Similar form for DFP update

How do we use LogDet?

同 ト イヨ ト イヨ ト

э

Kernel Learning

• The Matrix Nearness Problem:

min
$$D_{\ell d}(K, K_0)$$

 $(\mathbf{e}_i - \mathbf{e}_j)^T K(\mathbf{e}_i - \mathbf{e}_j) \leq u$ if $(i, j) \in S$ [similarity constraints]
 $(\mathbf{e}_i - \mathbf{e}_j)^T K(\mathbf{e}_i - \mathbf{e}_j) \geq \ell$ if $(i, j) \in D$ [dissimilarity constraints]
 $K \succeq 0$

э

伺 ト イヨト イヨト

Successive Projection-Correction Algorithm

- Algorithm: project successively onto each linear constraint followed by correction converges to globally optimal solution
- Each projection updates the Kernel matrix:

$$\begin{split} \min_{\mathcal{K}} & D_{\ell d}(\mathcal{K}, \mathcal{K}_t) \\ \text{s.t.} & (\mathbf{e}_i - \mathbf{e}_j)^{\mathsf{T}} \mathcal{K}(\mathbf{e}_i - \mathbf{e}_j) \leq u \end{split}$$

• Can be solved by rank-one update:

$$\begin{aligned} & \mathcal{K}_{t+1} = \mathcal{K}_t + \theta_t \mathcal{K}_t (\mathbf{e}_i - \mathbf{e}_j) (\mathbf{e}_i - \mathbf{e}_j)^T \mathcal{K}_t \\ \Rightarrow \mathcal{K}^* = \mathcal{K}_0 + \mathcal{K}_0 \Theta \mathcal{K}_0 \end{aligned}$$

- Advantages:
 - Automatic enforcement of positive semidefiniteness
 - Simple, closed-form projections (θ_t computable in closed form)
 - No eigenvalue/eigenvector calculation
 - Easy to incorporate slack for each constraint

Enhancements

- Impose structure on K:
 - Iow-rank
 - K_0 + low-rank
 - Each iteration costs only $O(r^2)$, where r is rank
- Extension to new data points:
 - Learned kernel can be shown to be of the form

$$\mathcal{K}(\mathbf{x},\mathbf{y}) = \mathcal{K}_0(\mathbf{x},\mathbf{y}) + \sum_i \sum_j heta_{ij} \mathcal{K}_0(\mathbf{x},\mathbf{x}_i) \mathcal{K}_0(\mathbf{y},\mathbf{x}_j)$$

- Fast Nearest-Neighbor Search
 - Can incorporate "locality-sensitive hashing"
- Online algorithms
 - Constraints are presented incrementally

Results

Inderjit S. Dhillon University of Texas at Austin The Log-Determinant Divergence and its Applications

▲御▶ ▲ 臣▶ ▲ 臣▶

æ

Application 1: Image Recognition

Data Set: Caltech 101

- Standard benchmark for multi-category image recognition
- 101 classes of images
- Wide variance in pose etc.
- Challenging data set

Experimental Setup

- 5, 10, 15 & 30 images per class in training set; rest in test set
- Performed 1-NN using original kernels and learned kernels



The Log-Determinant Divergence and its Applications

Results: Image Recognition



Inderjit S. Dhillon University of Texas at Austin The Log-Determinant Divergence and its Applications

Application 2: Automating Software Support



- Clarify: system to compare a user's program execution to existing executions in a database to automate software support
- Need appropriate notion of "distance" between program executions

伺 ト イ ヨ ト イ ヨ ト

Results: Clarify

- Representation: System collects program features during run-time
 - Function counts
 - Call-site counts
 - Counts of program paths
 - Program execution represented as a vector of counts
- Class labels: Program execution errors
- Nearest neighbor software support
 - Match program executions
 - Underlying distance measure should reflect this similarity

LaTeX Results



Conclusions

LogDet Divergence

- Previously used in Statistics & Optimization
- Has intriguing properties

Applications

- Leads to new matrix nearness problems
- Successive projection-correction algorithm
- Each projection can be computed efficiently

Challenges

- Faster solutions to matrix nearness problem interior point methods?
- Apply to find low-rank correlation matrices

References

"Learning Low-Rank Kernel Matrices", B. Kulis, M. A. Sustik, and I. S. Dhillon. International Conference on Machine Learning (ICML), pages 505-512, July 2006 (longer version submitted to journal).

"Information-Theoretic Metric Learning", J. Davis, B. Kulis, P. Jain, S. Sra, and I. S. Dhillon. International Conference on Machine Learning (ICML), pages 209-216, June 2007.

"Structured Metric Learning for High-Dimensional Problems", J. Davis and I. S. Dhillon. ACM International Conference on Knowledge Discovery and Data Mining(KDD), August 2008.

"Fast Image Search for Learned Metrics", P. Jain, B. Kulis, and K. Grauman. IEEE Conference on Computer Vision and Pattern Recognition (CVPR), June 2008.

"The Pyramid Match Kernel: Efficient Learning with Sets of Features", K. Grauman and T. Darrell. Journal of Machine Learning Research (JMLR), vol. 8, pages 725–760, Apr 2007.

"Matrix Nearness Problems using Bregman Divergences", I. S. Dhillon and J. A. Tropp. SIAM Journal on Matrix Analysis and its Applications, vol. 29, no. 4, pages 1120-1146, Dec 2007.

(日)

3