

Distance Functions and Metric Learning: Part 1

Michael Werman Ofir Pele



12/1/2011 version. Updated Version: http://www.cs.huji.ac.il/~ofirpele/DFML_ECCV2010_tutorial/

Rough Outline

- Bin-to-Bin Distances.
- Cross-Bin Distances:
 - Quadratic-Form (aka Mahalanobis) / Quadratic-Chi.
 - The Earth Mover's Distance.
- Perceptual Color Differences.
- Hands-On Code Example.

Distance ?



Metric

- Non-negativity: $D(P,Q) \ge 0$
- Identity of indiscernibles: D(P,Q) = 0 iff P = Q
- Symmetry: D(P,Q) = D(Q,P)
- Subadditivity (triangle inequality): $D(P,Q) \le D(P,K) + D(K,Q)$

Pseudo-Metric (aka Semi-Metric)

- Non-negativity: $D(P,Q) \ge 0$
- Property changed to: D(P,Q) = 0 (if) P = Q
- Symmetry: D(P,Q) = D(Q,P)
- Subadditivity (triangle inequality): $D(P,Q) \leq D(P,K) + D(K,Q)$

Metric

• Non-negativity: $D(P,Q) \ge 0$

• Identity of indiscernibles: D(P,Q) = 0 iff P = Q

An object is most similar to itself

Metric

Symmetry: D(P,Q) = D(Q,P)
Subadditivity (triangle inequality): D(P,Q) ≤ D(P,K) + D(K,Q)

Useful for many algorithms

Minkowski-Form Distances

 $rac{1}{p}$ $L_p(P,Q) = \sum |P_i - Q_i|^p$ i

Minkowski-Form Distances

 $L_1(P,Q) = \sum |P_i - Q_i|$

 $L_2(P,Q) = \sqrt{\sum_i (P_i - Q_i)^2}$

 $\overline{L_{\infty}(P,Q)} = \max_{i} |P_{i} - Q_{i}|$

Kullback-Leibler Divergence

$$KL(P,Q) = \sum_{i} P_i \log \frac{P_i}{Q_i}$$

- Information theoretic origin.
- Non symmetric.
- $Q_i = 0$?

Jensen-Shannon Divergence

$$JS(P,Q) = \frac{1}{2}KL(P,M) + \frac{1}{2}KL(Q,M)$$
$$M = \frac{1}{2}(P+Q)$$

$$JS(P,Q) = \frac{1}{2} \sum_{i} P_i \log \frac{2P_i}{P_i + Q_i} + \frac{1}{2} \sum_{i} Q_i \log \frac{2Q_i}{P_i + Q_i}$$

Jensen-Shannon Divergence

$$JS(P,Q) = \frac{1}{2}KL(P,M) + \frac{1}{2}KL(Q,M)$$
$$M = \frac{1}{2}(P+Q)$$

- Information Theoretic origin.
- Symmetric.
- \sqrt{JS} is a metric.

Jensen-Shannon Divergence

• Using Taylor extension and some algebra:

$$JS(P,Q) = \sum_{n=1}^{\infty} \frac{1}{2n(2n-1)} \sum_{i} \frac{(P_i - Q_i)^{2n}}{(P_i + Q_i)^{2n-1}}$$
$$= \frac{1}{2} \sum_{i} \frac{(P_i - Q_i)^2}{(P_i + Q_i)} +$$
$$\frac{1}{12} \sum_{i} \frac{(P_i - Q_i)^4}{(P_i + Q_i)^3} + \dots$$

$\chi^2\,$ Histogram Distance

 $\chi^{2}(P,Q) = \frac{1}{2} \sum_{i} \frac{(P_{i} - Q_{i})^{2}}{(P_{i} + Q_{i})}$

- Statistical origin.
- Experimentally results are very similar to JS.
- Reduces the effect of large bins.
- $\sqrt{\chi^2}$ is a metric

χ^2 Histogram Distance



χ^2 Histogram Distance



χ^2 Histogram Distance

 $\chi^{2}(P,Q) = \frac{1}{2} \sum_{i} \frac{(P_{i} - Q_{i})^{2}}{(P_{i} + Q_{i})}$

• Experimentally better than L_2 .

Bin-to-Bin Distances

• Bin-to-Bin distances such as L_1, L_2, χ^2 are sensitive to quantization:



Bin-to-Bin Distances

#bins robustness distinctiveness
#bins robustness distinctiveness



Bin-to-Bin Distances

#bins robustness distinctiveness
 #bins robustness distinctiveness

Can we achieve robustness and distinctiveness?

$$QF^{A}(P,Q) = \sqrt{(P-Q)^{T}A(P-Q)}$$

$$= \sqrt{\sum_{ij} (P_i - Q_i)(P_j - Q_j)A_{ij}}$$

- A_{ij} is the similarity between bin i and j.
- If *A* is the inverse of the covariance matrix, QF is called Mahalanobis distance.

$$QF^{A}(P,Q) = \sqrt{(P-Q)^{T}A(P-Q)}$$

$$= \sqrt{\sum_{ij} (P_i - Q_i)(P_j - Q_j)A_{ij}}$$

$$A = I$$

$$= \sqrt{\sum_{ij} (P_i - Q_i)^2} = L_2(P, Q)$$

$$QF^{A}(P,Q) = \sqrt{(P-Q)^{T}A(P-Q)}$$

- Does not reduce the effect of large bins.
- Alleviates the quantization problem.
- Linear time computation in # non zero A_{ij} .

$$QF^{A}(P,Q) = \sqrt{(P-Q)^{T}A(P-Q)}$$

• If A is positive-semidefinite then QF is a pseudo-metric.

Pseudo-Metric (aka Semi-Metric)

- Non-negativity: $D(P,Q) \ge 0$
- Property changed to: D(P,Q) = 0 (if) P = Q
- Symmetry: D(P,Q) = D(Q,P)
- Subadditivity (triangle inequality): $D(P,Q) \leq D(P,K) + D(K,Q)$

$$QF^{A}(P,Q) = \sqrt{(P-Q)^{T}A(P-Q)}$$

• If A is positive-definite then QF is a metric.

$$QF^{A}(P,Q) = \sqrt{(P-Q)^{T}A(P-Q)}$$

$$= \sqrt{(P-Q)^T W^T W (P-Q)}$$

 $=L_2(WP,WQ)$

$$QF^{A}(P,Q) = \sqrt{(P-Q)^{T}A(P-Q)}$$

$=L_2(WP,WQ)$

- We assume there is a linear transformation that makes bins independent
- There are cases where this is not true



$$QF^{A}(P,Q) = \sqrt{(P-Q)^{T}A(P-Q)}$$

• Converting distance to similarity (Hafner *et. al* 95):

$$A_{ij} = 1 - \frac{D_{ij}}{\max_{ij}(D_{ij})}$$

$$QF^{A}(P,Q) = \sqrt{(P-Q)^{T}A(P-Q)}$$

• Converting distance to similarity (Hafner *et. al* 95):

$$A_{ij} = e^{-\alpha \frac{D(i,j)}{\max_{ij}(D(i,j))}}$$

If α is large enough, A will be positive-definitive

$$QF^{A}(P,Q) = \sqrt{(P-Q)^{T}A(P-Q)}$$

• Learning the similarity matrix: part 2 of this tutorial.

$$QC_{m}^{A}(P,Q) = \sqrt{\sum_{ij} \frac{(P_{i} - Q_{i})(P_{j} - Q_{j})A_{ij}}{(\sum_{c}(P_{c} + Q_{c})A_{ci})^{m}(\sum_{c}(P_{c} + Q_{c})A_{cj})^{m}}}$$

- A_{ij} is the similarity between bin i and j.
- Generalizes QF and χ^2 .
- Reduces the effect of large bins.
- Alleviates the quantization problem.
- Linear time computation in # non zero A_{ij} .

$$QC_{m}^{A}(P,Q) = \sqrt{\sum_{ij} \frac{(P_{i} - Q_{i})(P_{j} - Q_{j})A_{ij}}{(\sum_{c} (P_{c} + Q_{c})A_{ci})^{m}(\sum_{c} (P_{c} + Q_{c})A_{cj})^{m}}}$$

is non-negative if A is positive-semidefinite.

- Symmetric.
- Triangle inequality unknown.
- If we define $\frac{0}{0} = 0$ and $0 \le m < 1$, QC is continuous.

Similarity-Matrix-Quantization-Invariant Property



Sparseness-Invariant





function dist= QuadraticChi(P,Q,A,m)

```
Z= (P+Q)*A;
% 1 can be any number as Z_i==0 iff D_i=0
Z(Z==0)= 1;
Z= Z.^m;
D= (P-Q)./Z;
% max is redundant if A is
% positive-semidefinite
dist= sqrt( max(D*A*D',0) );
```


Time Complexity: $O(\#A_{ij} \neq 0)$

function dist= QuadraticChi(P,Q,A,m)

```
Z= (P+Q)*A;
% 1 can be any number as Z_i==0 iff D_i=0
Z(Z==0)= 1;
Z= Z.^m;
D= (P-Q)./Z;
% max is redundant if A is
% positive-semidefinite
dist= sqrt( max(D*A*D',0) );
```



What about sparse (e.g. BoW) histograms ?

function dist= QuadraticChi(P,Q,A,m)

```
Z= (P+Q)*A;
% 1 can be any number as Z_i==0 iff D_i=0
Z(Z==0)= 1;
Z= Z.^m;
D= (P-Q)./Z;
% max is redundant if A is
% positive-semidefinite
dist= sqrt( max(D*A*D',0) );
```

The Quadratic-Chi Histogram Distance Code

What about sparse (e.g. BoW) histograms ?

P - Q = $0 \quad 0 \quad 0 \quad -3 \quad 0 \quad 0 \quad 4 \quad 5 \quad 0 \quad 0 \quad 0 \quad 0$ Time Complexity: O(SK)

The Quadratic-Chi Histogram Distance Code

What about sparse (e.g. BoW) histograms ?

P - Q =0 0 0 -3 0 0 0 4 5 0 0 0 0

Time Complexity: O(SK)

$\#(P \neq 0) + \#(Q \neq 0)$

The Quadratic-Chi Histogram Distance Code

What about sparse (e.g. BoW) histograms ?

P - Q =0 0 -3 0 0 -4 5 0 0 0 0 Time Complexity: O(SK)Average of non-zero entries in each row of A

 The Earth Mover's Distance is defined as the minimal cost that must be paid to transform one histogram into the other, where there is a "ground distance" between the basic features that are aggregated into the histogram.





 $EMD^{D}(P,Q) = \min_{F = \{F_{ij}\}} \sum_{i,j} F_{ij} D_{ij}$ $\sum F_{ij} = P_i \qquad \sum F_{ij} = Q_j$ s.t: $\sum F_{ij} = \sum P_i = \sum Q_j = 1$ $F_{ij} \ge 0$

 $EMD^{D}(P,Q) = \min_{F = \{F_{ij}\}} \frac{\sum_{i,j} F_{ij} D_{ij}}{\sum_{i} F_{ij}}$ $\sum F_{ij} \leq P_i \qquad \sum F_{ij} \leq Q_j$ s.t: $\sum F_{ij} = \min(\sum P_i, \sum Q_j)$ $F_{ij} \ge 0$

• Pele and Werman $08 - \widehat{EMD}$, a new EMD definition.

Definition:

 $\widehat{EMD}_C^D(P,Q) = \min_{\substack{F = \{F_{ij}\}\\i,j}} \sum_{i,j} F_{ij} D_{ij} +$ $\left|\sum_{i} P_{i} - \sum_{j} Q_{j}\right| \times C$ s.t: $\sum F_{ij} \leq P_i \qquad \sum F_{ij} \leq Q_j$ $\sum F_{ij} = \min(\sum P_i, \sum Q_j) \qquad F_{ij} \ge 0$ i, j

EMD





Demander



• When the total mass of two histograms is important.

$EMD(\square, \square) = EMD(\square, \square)$

• When the total mass of two histograms is important.

$$EMD(,) = EMD(,)$$



• When the difference in total mass between histograms is a distinctive cue.

$EMD(\square,\square) = EMD(\square,\square) = 0$

• When the difference in total mass between histograms is a distinctive cue.

 $EMD(\square,\square) = EMD(\square,\square) = 0$ $\widehat{EMD}(\square,\square) < \widehat{EMD}(\square,\square)$

- If ground distance is a metric:
- *EMD* is a metric only for normalized histograms.
 EMD is a metric for all histograms (C ≥ ¹/₂Δ).

 \widetilde{EMD} - a Natural Extension to L_1

• $EMD = L_1$ if: $D_{ij} = \begin{cases} 0 & \text{if } i = j \\ 2 & \text{otherwise} \end{cases}$ $\overline{C} = 1$

 $\widehat{EMD}_{C}^{D}(P,Q) = \min_{F = \{F_{ij}\}} \sum_{j \in i} F_{ij} D_{ij} +$ $\left|\sum_{i} P_{i} - \sum_{j} Q_{j}\right| \times C$

- General ground distance: $O(N^3 \log N)$ - Orlin 88
- L_1 normalized 1D histograms O(N)
 - Werman, Peleg and Rosenfeld 85

- L_1 normalized 1D cyclic histograms O(N)
 - Pele and Werman 08 (Werman, Peleg, Melter, and Kong 86)

240

300

360



• L_1 Manhattan grids $O(N^2 \log N(D + \log N))$ - Ling and Okada 07



• What about N-dimensional histograms with a cyclic dimensions?



• What about N-dimensional histograms with a cyclic dimensions?



• What about N-dimensional histograms with a cyclic dimensions?



- L_1 general histograms $O(N^2 \log^{2D-1} N)$
 - Gudmundsson, Klein, Knauer and Smid 07



- $\min(L_1, 2)$ 1D linear/cyclic $\{1, 2, \dots, \Delta\}$ O(N) histograms
 - Pele and Werman 08



- $\min(L_1, 2)$ 1D linear/cyclic $\{1, 2, \dots, \Delta\}$ O(N) histograms
 - Pele and Werman 08



t	*		
~	K		
	*	→	
1	1	>	7

- $\min(L_1, 2)$ general $\{1, 2, \dots, \Delta\}^D$ histograms *K* is the number of edges with cost $1O(N^2K\log(\frac{N}{K}))$
 - Pele and Werman 08

• Any thresholded distance $O(N^2 \log N(K + \log N))$ - Pele and Werman 09

> number of edges with cost different from the threshold

 $O(N^2 \log^2 N) \longleftarrow K = O(\log N)$

Thresholded Distances

• EMD with a thresholded ground distance is

not an approximation of EMD.

• It has better performance.

The Flow Network Transformation



Original Network

Simplified Network

The Flow Network Transformation



Original Network

Simplified Network

The Flow Network Transformation



Flowing the Monge sequence (if ground distance is a metric, zero-cost edges are a Monge sequence)
The Flow Network Transformation



The Flow Network Transformation



Combining Algorithms

- EMD algorithms can be combined.
- For example L_1 :



Combining Algorithms

- EMD algorithms can be combined.
- For example, thresholded L_1 :



• Charikar 02, Indyk and Thaper 03 – approximated EMD on $\{1, \ldots, \Delta\}^d$ by embedding it into the L_1 norm.



Time complexity: $O(TNd \log \Delta)$ Distortion (in expectation): $O(d \log \Delta)$

- Grauman and Darrell 05 Pyramid Match Kernel (PMK) same as Indyk and Thaper, replacing L_1 with histogram intersection.
- PMK approximates EMD with partial matching.
- PMK is a mercer kernel.
- Time complexity & distortion same as Indyk and Thaper (proved in Grauman and Darrell 07).

 Lazebnik, Schmid and Ponce 06 – used PMK in the spatial domain (SPM).







 Shirdhonkar and Jacobs 08 - approximated EMD using the sum of absolute values of the weighted wavelet coefficients of the difference histogram.



- Khot and Naor 06 any embedding of the EMD over the d-dimensional Hamming cube into L_1 must incur a distortion of $\Omega(d)$.
- Andoni, Indyk and Krauthgamer 08 for sets with cardinalities upper bounded by a parameter s the distortion reduces to $O(\log s \log d)$.
- Naor and Schechtman 07 any embedding of the EMD over $\{0,1,\ldots,\Delta\}^2$ must incur a distortion of $\Omega(\sqrt{\log\Delta})$.

Robust Distances

• Very high distances \Longrightarrow outliers \Longrightarrow same difference.



Robust Distances

• With colors, the natural choice.



Robust Distances

$\Delta E_{00}(blue, red) = 56$

ΔE_{00} (blue, yellow) = 102

Robust Distances - Exponent

- Usually a negative exponent is used:
 - Let d(a, b) be a distance measure between two features - a, b.
 - The negative exponent distance is:

$$d_e(a,b) = 1 - e^{\frac{-d(a,b)}{\sigma}}$$

Robust Distances - Exponent

 Exponent is used because (Ruzon and Tomasi 01): robust, smooth, monotonic, and a metric
Input is always discrete anyway ...

Robust Distances - Thresholded

- Let d(a, b) be a distance measure between two features - a, b.
- The thresholded distance with a threshold of t > 0 is:

$$d_t(a,b) = \min(d(a,b),t).$$

- Thresholded metrics are also metrics (Pele and Werman ICCV 2009).
- Better results.
- Pele and Werman ICCV 2009 algorithm computes EMD with thresholded ground distances much faster.
- Thresholded distance corresponds to sparse similarities matrix -> faster QC / QF computation.

$$A_{ij} = 1 - \frac{D_{ij}}{\max_{ij}(D_{ij})}$$

- Thresholded vs. exponent:
 - Fast computation of cross-bin distances with a thresholded ground distance.
 - Exponent changes small distances can be a problem (e.g. color differences).

Color distance should be thresholded (robust).





The ground distance between two SIFT bins (x_i, y_i, o_i) and (x_j, y_j, o_j) :

$$d_{R} = ||(x_{i}, y_{i}) - (x_{j}, y_{j})||_{2} + \min(|o_{i} - o_{j}|, M - |o_{i} - o_{j}|)$$
$$d_{T} = \min(d_{R}, T)$$

A Ground Distance for Color Image

The ground distances between two LAB image bins $(x_i, y_i, L_i, a_i, b_i)$ and $(x_j, y_j, L_j, a_j, b_j)$ we use are:

$$dc_T = \min((||(x_i, y_i) - (x_j, y_j)||_2) + \Delta_{00}((L_i, a_i, b_i), (L_j, a_j, b_j)), T)$$

 Euclidean distance on L*a*b* space is widely considered as perceptual uniform.





• ΔE_{00} on L*a*b* space is better.

Luo, Cui and Rigg 01. Sharma, Wu and Dalal 05.



• ΔE_{00} on L*a*b* space is better.



• ΔE_{00} on L*a*b* space is better.



- ΔE_{00} on L*a*b* space is better.
- But still has major problems.

$\Delta E_{00}(blue, red) = 56$ $\Delta E_{00}(blue, red) = 102$

Color distance should be thresholded (robust).

Color distance should be saturated (robust).





Perceptual Color Descriptors

Perceptual Color Descriptors

• 11 basic color terms. Berlin and Kay 69.



Perceptual Color Descriptors

• 11 basic color terms. Berlin and Kay 69.


• 11 basic color terms. Berlin and Kay 69.



 Image copyright by Eric Rolph. Taken from: http://upload.wikimedia.org/wikipedia/commons/5/5c/Double-alaskan-rainbow.jpg

• 11 basic color terms. Berlin and Kay 69.



 How to give each pixel an "11-colors" description ?

• Learning Color Names from Real-World Images, J. van de Weijer, C. Schmid, J. Verbeek CVPR 2007.





• Learning Color Names from Real-World Images, J. van de Weijer, C. Schmid, J. Verbeek CVPR 2007.



















• Learning Color Names from Real-World Images, J. van de Weijer, C. Schmid, J. Verbeek CVPR 2007.



- Learning Color Names from Real-World Images, J. van de Weijer, C. Schmid, J. Verbeek CVPR 2007.
- For each color returns a probability distribution over the 11 basic colors.



- Applying Color Names to Image Description, J. van de Weijer, C. Schmid ICIP 2007.
- Outperformed state of the art color descriptors.



- Using illumination invariants black, gray and white are the same.
- "Too much invariance" happens in other cases

(Local features and kernels for classification of texture and object categories: An in-depth study - Zhang, Marszalek, Lazebnik and Schmid. IJCV 2007, Learning the discriminative power-invariance trade-off - Varma and Ray. ICCV 2007).

• To conclude: Don't solve imaginary problems.

- This method is still not perfect.
- 11 color vector for purple(255,0,255) is:



• In real world images there are no such over-saturated colors.

Open Questions

- EMD variant that reduces the effect of large bins.
- Learning the ground distance for EMD.
- Learning the similarity matrix and normalization factor for QC.

Hands-On Code Example

http://www.cs.huji.ac.il/~ofirpele/FastEMD/code/

http://www.cs.huji.ac.il/~ofirpele/QC/code/

Tutorial:

http://www.cs.huji.ac.il/~ofirpele/DFML_ECCV2010_tutorial/

