

Notion of Distance

Metric Distance

Binary Vector Distances

Tangent Distance

Distance Measures

- Many pattern recognition/data mining techniques are based on similarity measures between objects
 - e.g., nearest-neighbor classification,
 - cluster analysis,
 - multi-dimensional scaling
- $s(i,j)$: similarity, $d(i,j)$: dissimilarity
- Possible transformations:
 - $d(i,j) = 1 - s(i,j)$ or
 - $d(i,j) = \sqrt{2 * (1 - s(i,j))}$
- Proximity is a general term to indicate similarity and dissimilarity
- Distance is used to indicate dissimilarity

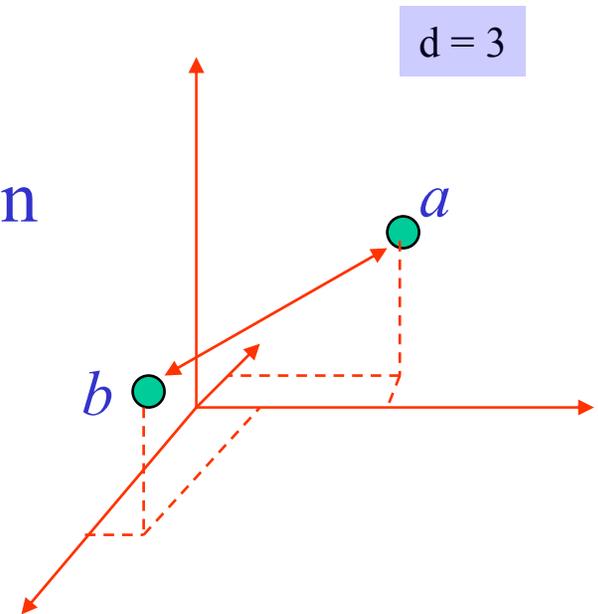
Euclidean distance

Nearest neighbor classifier relies on a distance function between patterns

The Euclidean formula for distance in d dimensions is

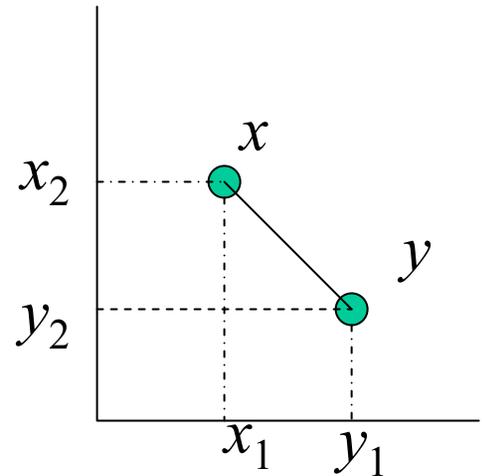
$$D(\mathbf{a}, \mathbf{b}) = \left(\sum_{k=1}^d (a_k - b_k)^2 \right)^{1/2}$$

Notion of a metric is far more general



Euclidean Distance between Vectors

$$d_E(x, y) = \left(\sum_{k=1}^p (x_k - y_k)^2 \right)^{1/2}$$



- Euclidean distance assumes variables are **commensurate**
- E.g., each variable a measure of length
- If one were weight and other was length there is no obvious choice of units
- Altering units would change which variables are important

Definition of a Metric

- Notion of a metric is more general than Euclidean distance
- Properties of a metric

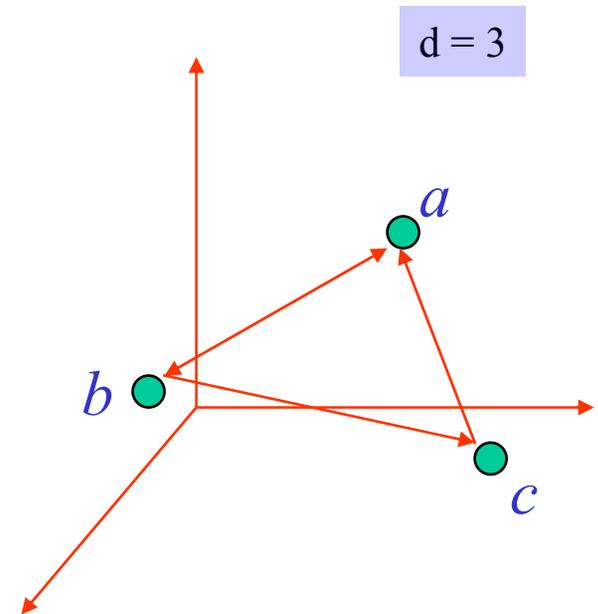
For all vectors **a**, **b** and **c**

nonnegativity: $D(\mathbf{a}, \mathbf{b}) \geq 0$

reflexivity: $D(\mathbf{a}, \mathbf{b}) = 0$ if and only if $\mathbf{a} = \mathbf{b}$

symmetry: $D(\mathbf{a}, \mathbf{b}) = D(\mathbf{b}, \mathbf{a})$

triangle inequality: $D(\mathbf{a}, \mathbf{b}) + D(\mathbf{b}, \mathbf{c}) \geq D(\mathbf{a}, \mathbf{c})$.

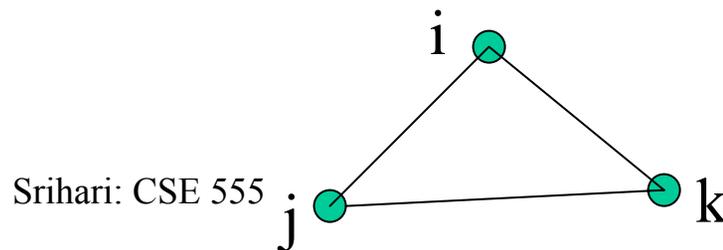


Euclidean distance is a metric

Metric Properties

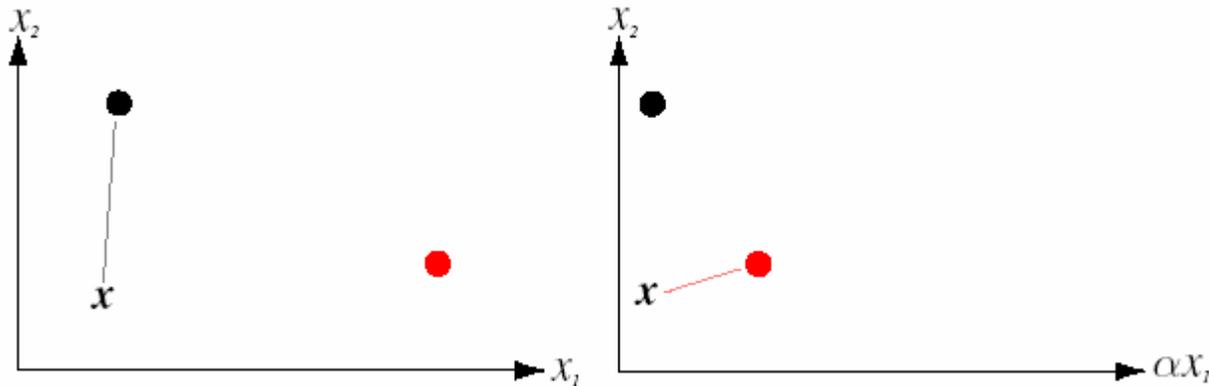
A **metric** is a dissimilarity (distance) measure that satisfies the following properties:

1. $d(i,j) \geq 0$ Positivity
2. $d(i,j) = d(j,i)$ Commutativity
3. $d(i,j) \leq d(i,k) + d(k,j)$ Triangle Inequality



Scale changes

can affect nearest neighbor classifiers based on Euclidean distance



Point x is closest to black

Point x is closest to Red

Rescaling the data to equalize ranges is equivalent to changing the metric in the original space

Standardizing the Data

when variables are not commensurate

- Divide each variable by its standard deviation
- Standard deviation for the k^{th} variable is

where
$$\sigma_k = \left(\frac{1}{n} \sum_{i=1}^n (x_k(i) - \mu_k)^2 \right)^{\frac{1}{2}}$$

Updated value that removes the effect of scale:

$$\mu_k = \frac{1}{n} \sum_{i=1}^n x_k(i)$$

$$x_k' = \frac{x_k}{\sigma_k}$$

Weighted Euclidean Distance

- If we know relative importance of variables

$$d_{WE}(i, j) = \left(\sum_{k=1}^p w_k ((x_k(i) - x_k(j))^2) \right)^{\frac{1}{2}}$$

Distance between Variables

- Similarities between cups



- Suppose we measure cup-height 100 times and diameter only once
 - height will dominate although 99 of the height measurements are not contributing anything
- They are very highly correlated
- To eliminate redundancy we need a *data-driven method*
 - approach is to not only to standardize data in each direction but also to use covariance between variables

Sample Covariance between variables X and Y

$$\text{Cov}(X, Y) = \frac{1}{n} \sum_{i=1}^n \left(x(i) - \bar{x} \right) \left(y(i) - \bar{y} \right)$$

Sample
means

- It is a scalar value that measures how X and Y vary together
- Obtained by multiplying for each sample its mean-centered value of x with mean-centered value of y and then adding over all samples
- Large **positive value** if large values of X tend to be associated with large values of Y and small values of X with small values of Y
- Large **negative value** if large values of X tend to be associated with small values of Y
- With p variables can construct a p x p matrix of covariances. Such a covariance matrix is symmetric.

Relationship between Covariance Matrix and Data Matrix

- Let $X = n \times p$ data matrix
- Rows of X are the data vectors $x(i)$
- Definition of covariance:

$$Cov(i, j) = \frac{1}{n} \sum_{k=1}^n \left(x_k(i) - \bar{x} \right) \left(y_k(i) - \bar{y} \right)$$

- If values of X are mean-centered (i.e., value of each variable is relative to the sample mean of that variable) then $V = X^T X$ is the $p \times p$ covariance matrix

Correlation Coefficient

Value of Covariance is **dependent** upon ranges of X and Y

Dependency is removed by dividing values of X by their standard deviation and values of Y by their standard deviation

$$\rho(X, Y) = \frac{\sum_{i=1}^n (x(i) - \bar{x})(y(i) - \bar{y})}{\sigma_x \sigma_y}$$

With p variables, can form a p x p **correlation matrix**

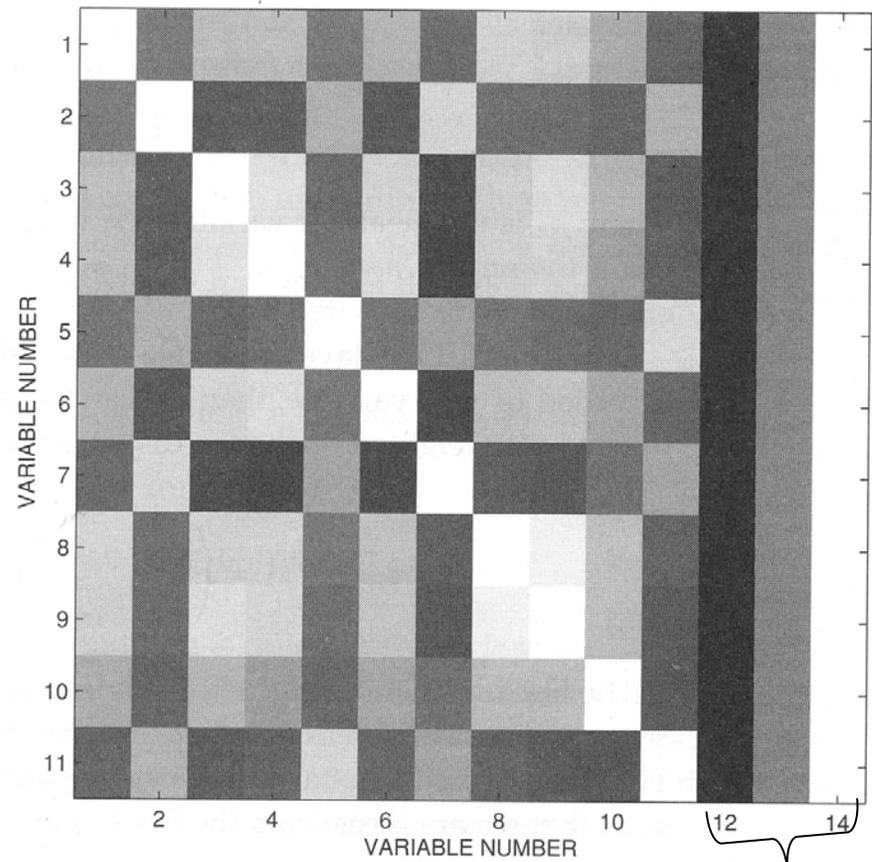
Correlation Matrix

(Housing related variables
across city suburbs)

Variables 3 and 4 are highly negatively
correlated with Variable 2

Variable 5 is positively correlated with
Variable 11

Variables 8 and 9 are highly correlated



Reference for -1, 0,+1

Figure 2.1 A sample correlation matrix plotted as a pixel image. White corresponds to +1 and black to -1. The three rightmost columns contain values of -1, 0, and +1 (respectively) to provide a reference for pixel intensities. The remaining 11 × 11 pixels represent the 11 × 11 correlation matrix. The data come from a well-known data set in the regression research literature, in which each data vector is a suburb of Boston and each variable represents a certain general characteristic of a suburb. The variable names are (1) per-capita crime rate, (2) proportion of area zoned for large residential lots, (3) proportion of non-retail business acres, (4) nitric oxide concentration, (5) average number of rooms perdwelling, (6) proportion of pre-1940 homes, (7) distance to retail centers index, (8) accessibility to highways index, (9) property tax rate, (10) pupil-to-teacher ratio, and (11) median value of owner-occupied homes.

Generalizing Euclidean Distance

Minkowski or L_λ metric $\left(\sum_{k=1}^p (x_k(i) - x_k(j))^\lambda \right)^{\frac{1}{\lambda}}$

- $\lambda = 2$ gives the Euclidean metric
- $\lambda = 1$ gives the Manhattan or City-block metric

$$\sum_{k=1}^p |x_k(i) - x_k(j)|$$

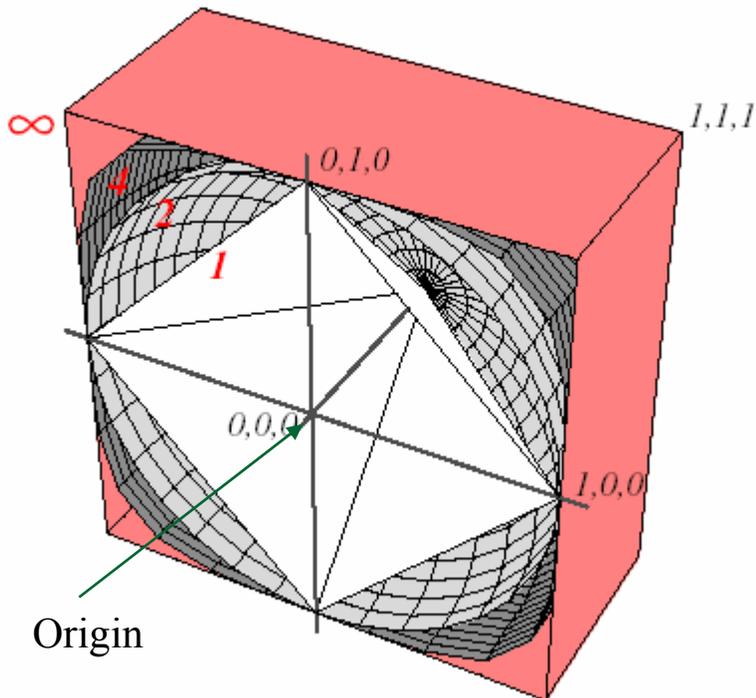
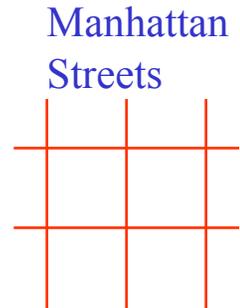
- $\lambda = \text{infinity}$ yields

$$\max_k |x_k(i) - x_k(j)|$$

Minkowski Metric

$$L_k(\mathbf{a}, \mathbf{b}) = \left(\sum_{i=1}^d |a_i - b_i|^k \right)^{1/k}$$

- The L_k norm
- L_1 norm is the Manhattan (city block) distance
- L_2 norm is the Euclidean distance



Each colored surface consists of points of distance 1.0 from the origin. Using different values for k in the Minkowski metric (k is in red).

Vector Space Representation of Documents

Document-Term Matrix

	t1	t2	t3	t4	t5	t6
D1	24	21	9	0	0	3
D2	32	10	5	0	3	0
D3	12	16	5	0	0	0
D4	6	7	2	0	0	0
D5	43	31	20	0	3	0
D6	2	0	0	18	7	16
D7	0	0	1	32	12	0
D8	3	0	0	22	4	2
D9	1	0	0	34	27	25
D10	6	0	0	17	4	23

Terms used
In "Database"
Related documents

t1	database
t2	SQL
t3	index
t4	regression
t5	likelihood
t6	linear

"Regression" Terms

d_{ij} represents number of times that term appears in that document

Cosine Distance between Document Vectors

Document-Term Matrix

	t1	t2	t3	t4	t5	t6
D1	24	21	9	0	0	3
D2	32	10	5	0	3	0
D3	12	16	5	0	0	0
D4	6	7	2	0	0	0
D5	43	31	20	0	3	0
D6	2	0	0	18	7	16
D7	0	0	1	32	12	0
D8	3	0	0	22	4	2
D9	1	0	0	34	27	25
D10	6	0	0	17	4	23

Cosine Distance

$$d_c(D_i, D_j) = \frac{\sum_{k=1}^T d_{ik} d_{jk}}{\sqrt{\sum_{k=1}^T d_{ik}^2 \sum_{k=1}^T d_{jk}^2}}$$

Cosine of the angle between two vectors

Equivalent to their inner product after each has been normalized to have unit length

Higher values for more similar vectors

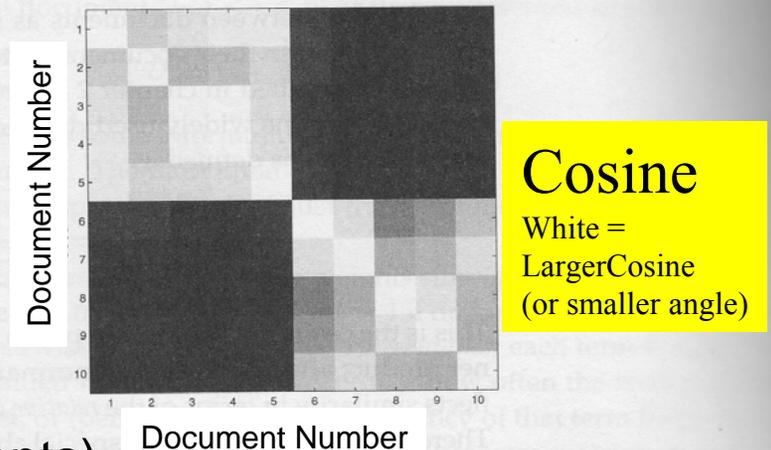
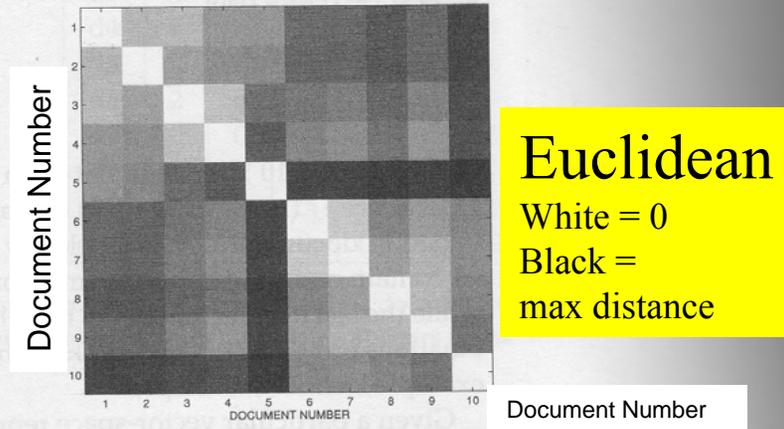
Reflects similarity in terms of relative distributions of components

Cosine is not influenced by one document being small compared to the other (as in the case of Euclidean)

Euclidean vs Cosine Distance

Document-Term Matrix

	t1	t2	t3	t4	t5	t6	
Database	D1	24	21	9	0	0	3
	D2	32	10	5	0	3	0
	D3	12	16	5	0	0	0
	D4	6	7	2	0	0	0
	D5	43	31	20	0	3	0
Regression	D6	2	0	0	18	7	16
	D7	0	0	1	32	12	0
	D8	3	0	0	22	4	2
	D9	1	0	0	34	27	25
	D10	6	0	0	17	4	23



Both show two clusters of light sub-blocks
(*database* documents and *regression* documents)

Euclidean: 3,4 closer to 6-9 than to 5 since 3,4, 6-9 are closer to origin than 5

Cosine emphasizes relative contributions of individual terms

Binary-Valued Feature Vectors

- Binary vectors are widely used in pattern recognition and information retrieval
- Several distance measures defined some of which are non-metric

Distance Measures for Binary Data

- Most obvious measure is *Hamming Distance* normalized by number of bits

$$\frac{S_{11} + S_{00}}{S_{11} + S_{10} + S_{01} + S_{00}}$$

- If we don't care about irrelevant properties had by neither object we have *Jaccard Coefficient*

$$\frac{S_{11}}{S_{11} + S_{10} + S_{01}}$$

- *Dice Coefficient* extends this argument. If 00 matches are irrelevant then 10 and 01 matches should have half relevance

Tanimoto Metric

- Used in Taxonomy

$$D_{Tanimoto}(S_1, S_2) = \frac{n_1 + n_2 - 2n_{12}}{n_1 + n_2 - n_{12}},$$

- n_1 and n_2 are the no of elements in S_1 and S_2
- n_{12} is the no that is in both sets
- Useful when the features are same or different (binary valued)
 - 0000 1111 0000 vs 1000 1101 0000

Similarity/Dissimilarity Measures for Binary Vectors

Measure	$S(X, Y)$	$D(X, Y)$
Jaccard-Needham	$\frac{S_{11}}{S_{11} + S_{10} + S_{01}}$	$\frac{S_{10} + S_{01}}{S_{11} + S_{10} + S_{01}}$
Dice	$\frac{S_{11}}{2S_{11} + S_{10} + S_{01}}$	$\frac{S_{10} + S_{01}}{2S_{11} + S_{10} + S_{01}}$
Correlation	$\frac{S_{11}S_{00} - S_{10}S_{01}}{\sigma}$	$\frac{1}{2} - \frac{S_{11}S_{00} - S_{10}S_{01}}{2\sigma}$
Yule	$\frac{S_{11}S_{00} - S_{10}S_{01}}{S_{11}S_{00} + S_{10}S_{01}}$	$\frac{S_{10}S_{01}}{S_{11}S_{00} + S_{10}S_{01}}$
* Russell-Rao	$\frac{S_{11}}{N}$	$\frac{N - S_{11}}{N}$
* Sokal-Michener	$\frac{S_{11} + S_{00}}{N}$	$1 - \frac{S_{11} + S_{00}}{N}$
Rogers-Tanmoto	$\frac{S_{11} + S_{00}}{S_{11} + S_{00} + 2S_{10} + 2S_{01}}$	$\frac{2S_{10} + 2S_{01}}{S_{11} + S_{00} + 2S_{10} + 2S_{01}}$
* Rogers-Tanmoto-a	$\frac{S_{11} + S_{00}}{S_{11} + S_{00} + 2S_{10} + 2S_{01}}$	$\frac{2(N - S_{11} - S_{00})}{2N - S_{11} - S_{00}}$
Kulzinsky	$\frac{S_{11}}{(S_{10} + S_{01})}$	$\frac{S_{10} + S_{01} - S_{11} + N}{(S_{10} + S_{01} + N)}$

where $\sigma = ((S_{10} + S_{11})(S_{01} + S_{00})(S_{11} + S_{01})(S_{00} + S_{10}))^{1/2}$

$S_{00}(X, Y) = \bar{X} \cdot \bar{Y}$, $S_{01}(X, Y) = \bar{X} \cdot Y$ $S_{10}(X, Y) = X \cdot \bar{Y}$, $S_{11}(X, Y) = X \cdot Y$

Weighted Dissimilarity Measures for Binary Vectors

- Unequal importance to '0' matches and '1' matches
- Multiply S_{00} with β ($[0,1]$)
- Examples:

$$D^{sm}(X, Y) = 1 - \frac{S_{11} + \beta \cdot S_{00}}{N}$$

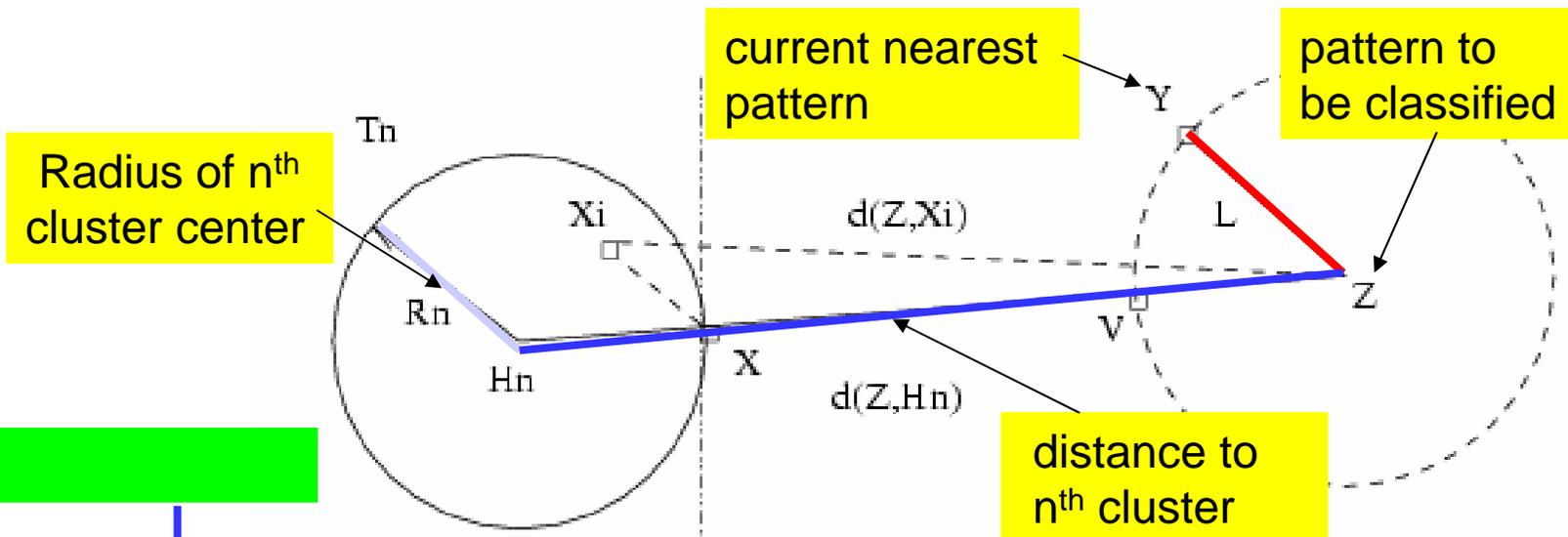
$$D^{rta}(X, Y) = \frac{2(N - S_{11} - \beta \cdot S_{00})}{2N - S_{11} - \beta \cdot S_{00}}$$

k-nn digit recognition accuracy

- Sokal-Michener ($\beta = 0.35$) → 99.38%
- Rogers-Tanimoto
Correlation
Jaccard-Needham
Dice } → 99.36%
- Yule → 99.28%
- Kulzinsky → 99.03%
- Russell-Rao → 97.56%

Speed Issues

k-nn Reduction Rules



Rule 1

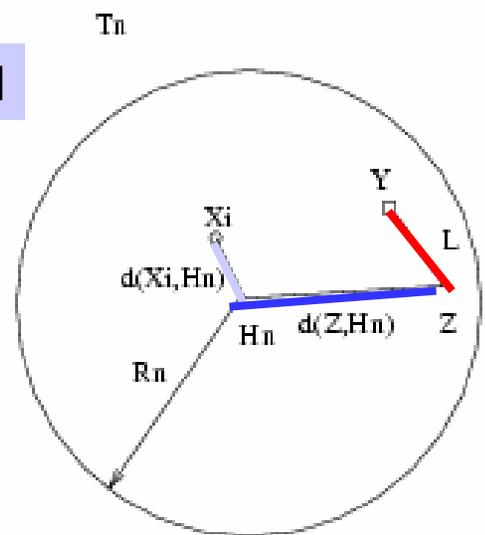
$$|L| + |R_n| < |d(Z, H_n)|$$

then no member of T_n need be considered

Rule 2

$$|L| + |d(X_i, H_n)| < |d(Z, H_n)|$$

then X_i cannot be the nearest neighbor



(a)

Tri-Edge Inequality (TEI)

N-dimensional binary space Ω

Θ is a subspace of Ω

$D(\cdot, \cdot)$ is a dissimilarity measure

TEI for $D(\cdot, \cdot)$ over Θ means that the sum of any two dissimilarities (edges) is equal to or bigger than any single dissimilarity:

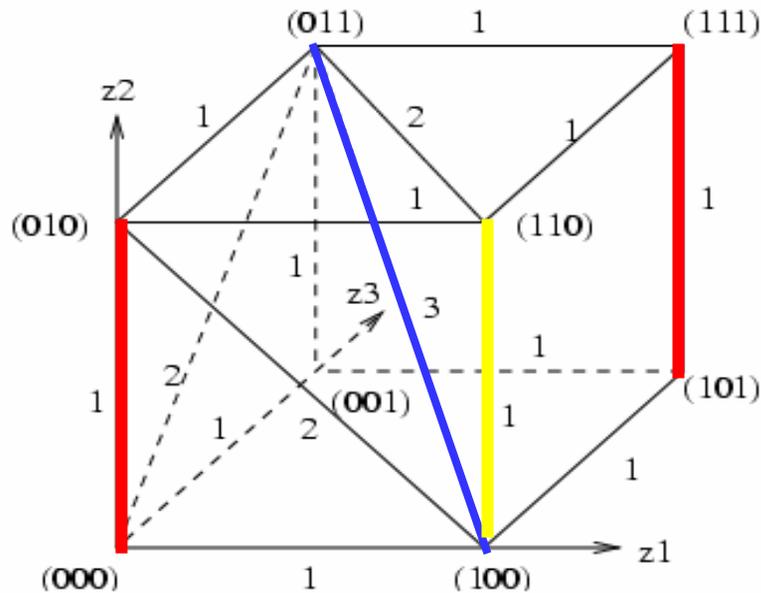
$$\forall X, Y, U, V, P, Q \in \Theta, \text{ and } X \neq Y \text{ and } U \neq V,$$

$$D(X, Y) + D(U, V) \geq D(P, Q)$$

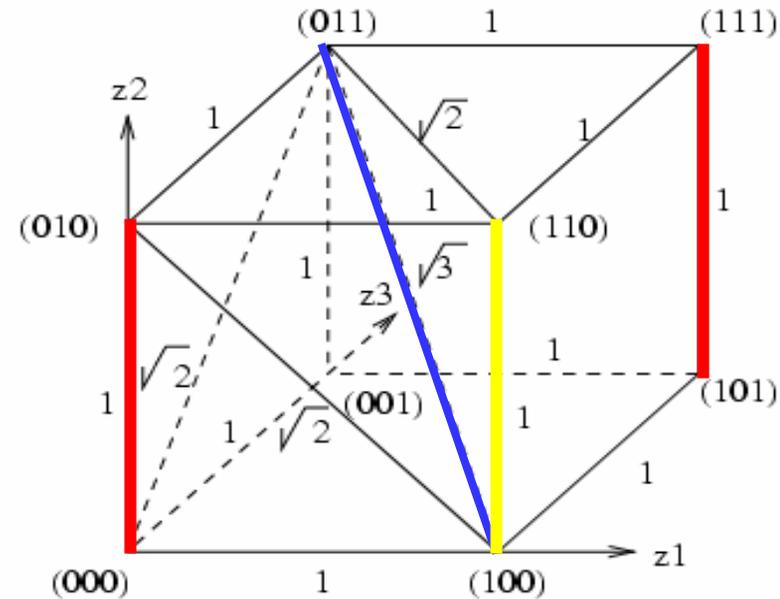
Note: here three edges, $D(X, Y)$, $D(U, V)$ and $D(P, Q)$, may not originate from a triangle, thus TEI is different from the Tri-Angle Inequality.

Tri Edge Inequality in 3-d Binary Space

Hamming Distance



Euclidean Distance



TEI holds for yellow edge, *not* for blue:

$$\begin{array}{l}
 \color{red}{|} (1) + \color{red}{|} (1) > \color{yellow}{|} (1) \\
 \color{blue}{|} (3) < \color{red}{|} (1) + \color{red}{|} (1)
 \end{array}$$

TEI holds (for all edges):
 $\min(D(,))=1$
 and $\max(D(,)) = \sqrt{3}$, thus,

$$2 * \min(D(,)) > \max(D(,))$$

$$2 * 1 > \sqrt{3}$$

TEI Property with Binary Vector Dissimilarity Measures

Measure	metric	Tri-Edge Inequality
Jaccard-Needham	unknown	no
Dice	no	no
Correlation	no	no
Yule	no	no
Russell-Rao	no	$M \leq N/2$
Sokal-Michener	$\beta = 1$	$\beta \leq \frac{1}{1 + \frac{2\lambda}{N-2M}}$
Rogers-Tanmoto	$\beta = 1$	no
Rogers-Tanmoto-a	$\beta = 1$	yes
Kulzinsky	unknown	no

Test for TEI

Definition:

$\Omega = \{X \mid X \text{ is an } N\text{-dimensional binary column vector}\}$

$\Theta \subseteq \Omega$ where

$$\Theta(\lambda, M) = \{X \mid X \in \Omega, y \in \Theta, \|X\| \leq \lambda, X^t \cdot Y \leq M\}$$

The Test:

Let $D_{min}(\cdot, \cdot) = \min_{X, Y \in \Theta} D(X, Y)$ and $D_{max}(\cdot, \cdot) = \max_{X, Y \in \Theta} D(X, Y)$. Then,

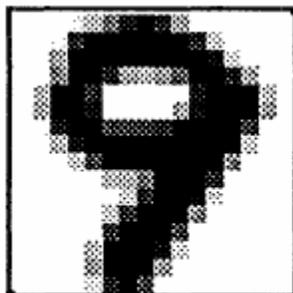
the tri-edge inequality can be verified through checking $2D_{min}(\cdot, \cdot) \geq D_{max}(\cdot, \cdot)$,

which is a much stronger condition.

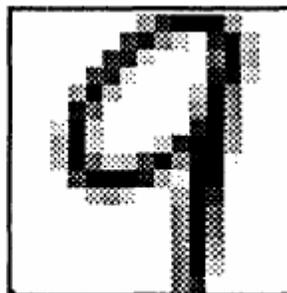
Summary of Binary Vector measures

- TEI is a property of the dissimilarity measure for binary vector classification
- It holds with three binary vector dissimilarity measures
- TEI test defined and experimentally verified
- Reduction rules of fast k-nn algorithms are inapplicable to dissimilarity measures violating TEI

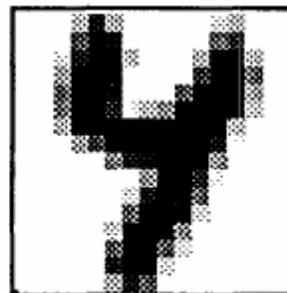
Tangent Distance



Pattern to
be classified



Prototype A



Prototype B

Which is more similar to test pattern, A or B?

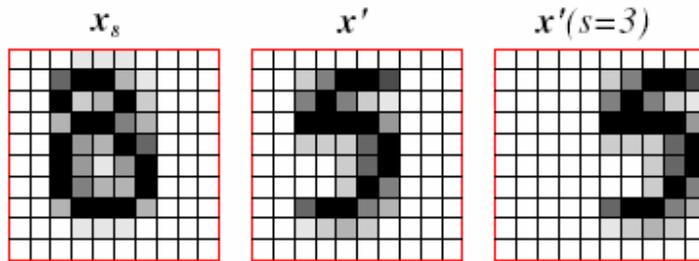
Euclidean distance: sum of squares of pixel to pixel difference

According to Euclidean distance, Prototype B is more similar

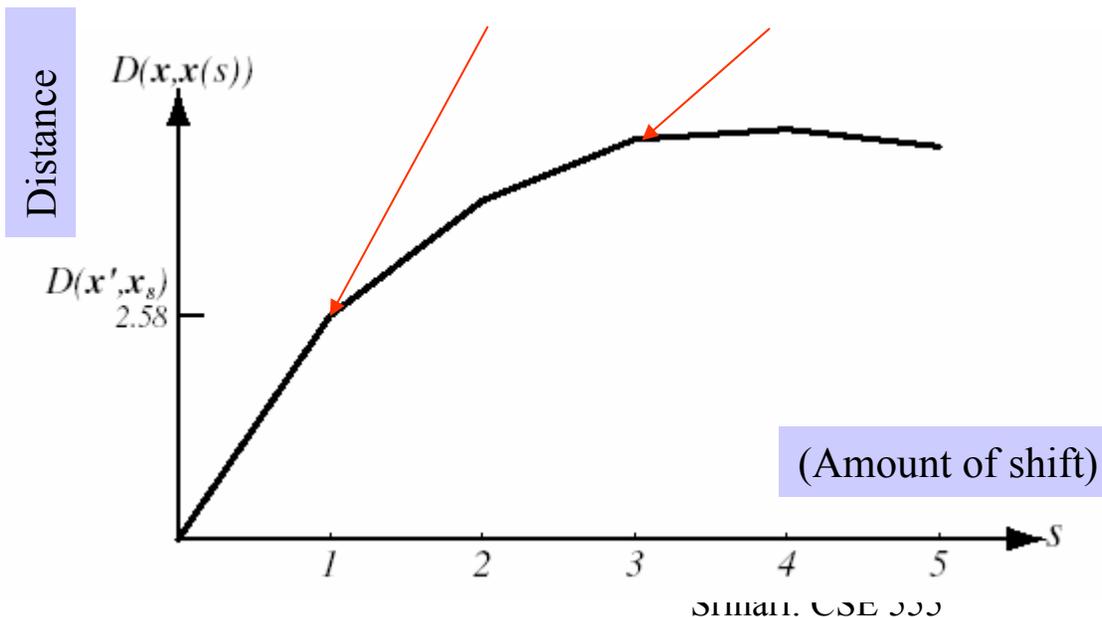
Even though prototype A is more similar once
it has been rotated and thickened

Effect of Translation on Euclidean Distance

- Pattern vector: 10x10 (d=100) grey values



When shift is large
Distance between
5's is more than
Distance between
5 and 8

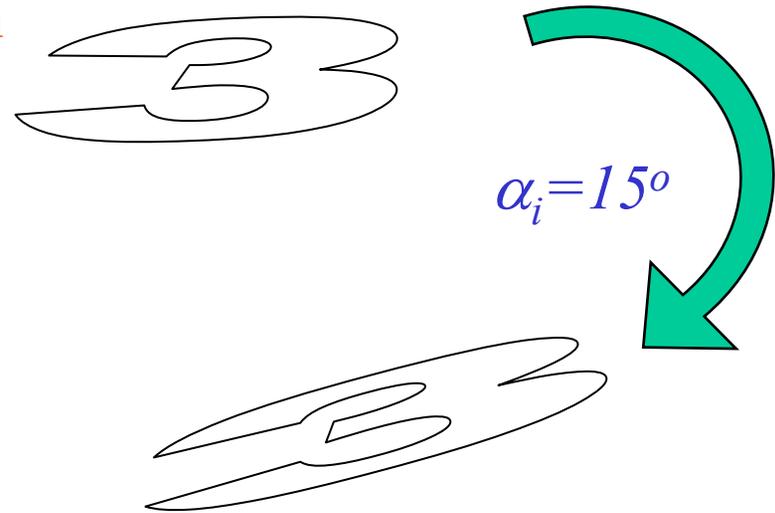


Euclidean
Distance is not
Translation
invariant

Transformations

- Suppose there are r possible transformations:

- Horizontal (x) translation
- Vertical (y) translation
- Shear
- Rotation
- Scale
- Thinning



- For each prototype \mathbf{x}' perform each transformation $F_i(\mathbf{x}'; \alpha_i)$
 - E.g., $F_i(\mathbf{x}'; \alpha_i)$ represents \mathbf{x}' rotated by α_i

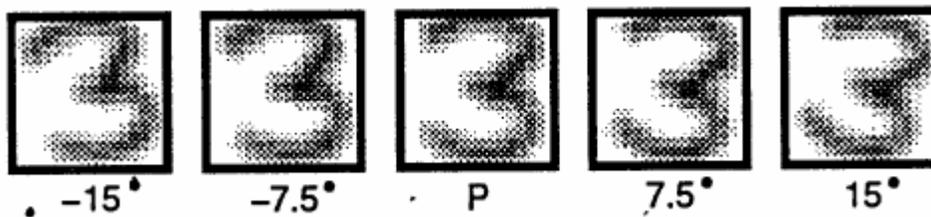
Transformation that depends on one parameter α (rotation angle)

Pattern P is transformed with a transformation s that depends on one parameter α

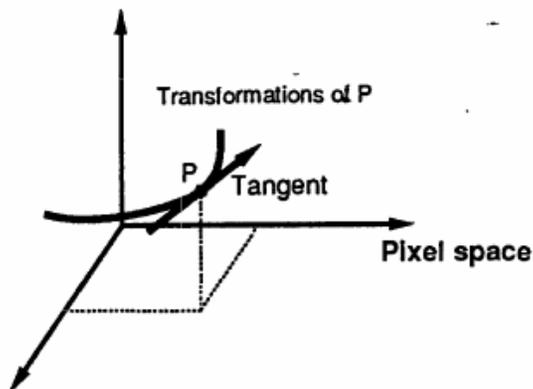
Set of all transformed patterns $S_P = \{x \mid \exists \alpha \text{ such that } x = s(\alpha, P)\}$

is a one - dimensional curve in vector space of inputs

Small rotations
of original
digit 3

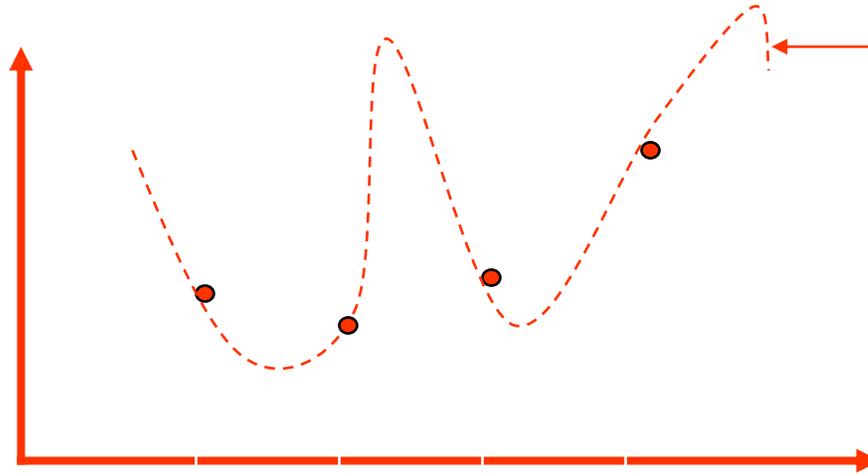


Effect of the
rotations in
pixel space
(if there were only
3 pixels)

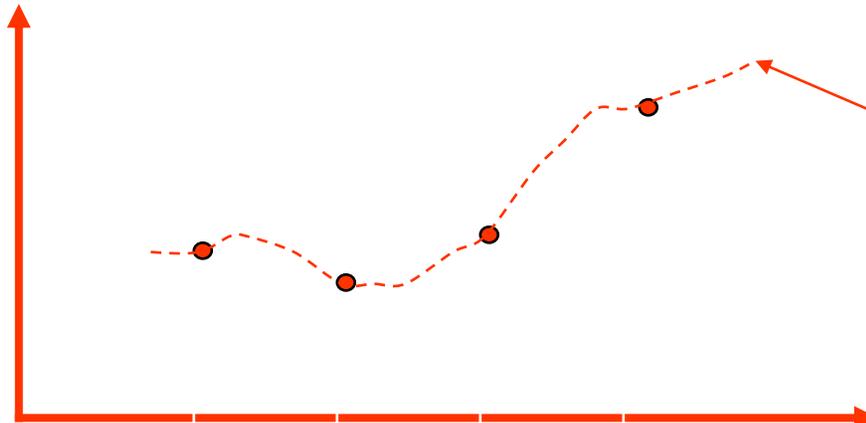


When there are n parameters α_i (rotation, scaling, etc)
 S_P is a manifold of at most n dimensions

Learning from a limited set of samples

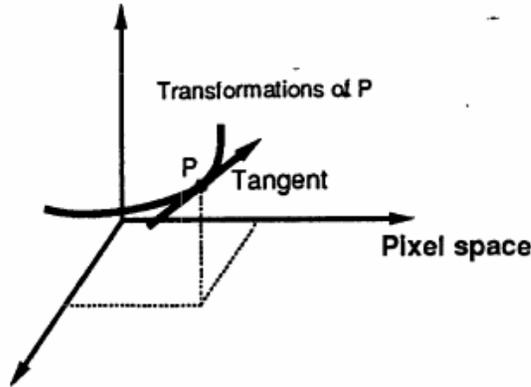


Fitted curve with only constraint that it goes through the samples: a poor fit



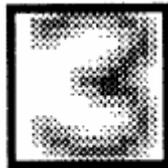
Fitted curve with constraint that it goes through the samples and its Derivatives: a better fit

Tangent Vector



Small Transformations of P can be obtained by adding to P a linear combination of vectors that span the tangent plane (tangent vectors)

Tangent vectors for a transformation s can be easily computed by finite difference, by evaluating $\delta s(P, \alpha) / \delta \alpha$



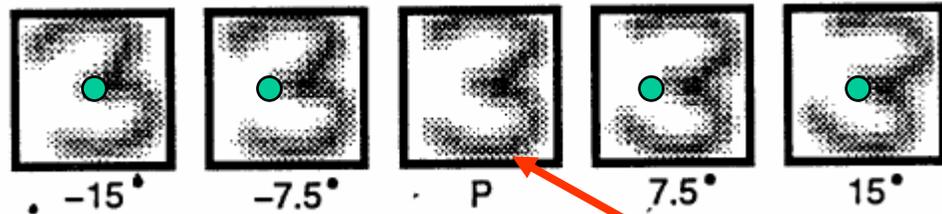
P



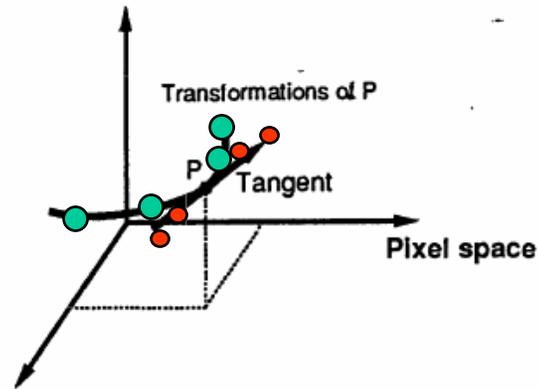
T. V.

Transformation that depends on one parameter α (angle of rotation)

Small rotations of original digit 3

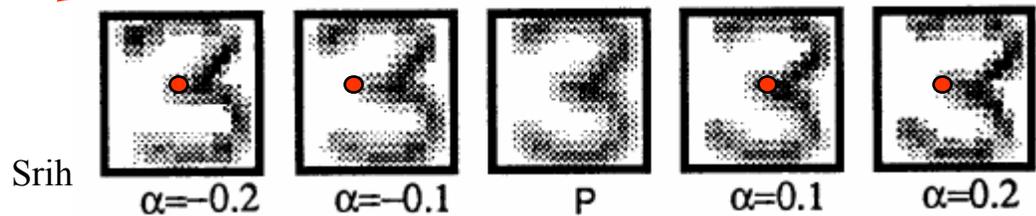
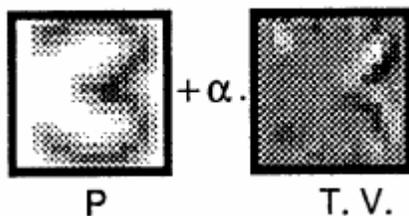


Effect of the rotations in pixel space (if there were only 3 pixels)



Same Original Image

Images obtained by moving along the tangent to the transformation curve for the same original digitized image P by adding various amounts (α) of the tangent vector TV



Tangent Vector

- Construct a Tangent vector TV_i for each transformation:

$$\mathbf{TV}_i = \mathcal{F}_i(\mathbf{x}'; \alpha_i) - \mathbf{x}'.$$

- Assume \mathbf{x}' is represented by a d -dimensional vector
- For each prototype \mathbf{x}' construct an $r \times d$ matrix \mathbf{T} consisting of the tangent vectors at \mathbf{x}'
- If there are 2 transformations and a 256-element vector (16 x 16 pixels) is used then there are *two* tangent vector images
- The tangent vectors are computed at training time

Algorithm for computing tangent vectors

- $T(a,u)$ denotes the input u rotated by α
- By definition the tangent to the transformation is given by

$$T_u = \left. \frac{\partial t(\alpha, u)}{\partial \alpha} \right|_{\alpha=0}$$

$$T_u \approx \frac{t(\varepsilon, u) - t(0, u)}{\varepsilon}$$

- This approach is not accurate, also $t(\varepsilon, u)$ for 2-d images requires a 2-d interpolation which is expensive

Tangent Vector for 2-D images

- Let $u(x,y)$ be an image
- T be a transformation that displaces the coordinates X and Y of each pixel by D_x and D_y and smoothes it by convolving it with a Gaussian g

$$t(u, \alpha) \approx u * g + \alpha \left(D_x \frac{\partial(u * g)}{\partial x} + D_y \frac{\partial(u * g)}{\partial y} \right)$$

$$t(u, \alpha) \approx u * g + \alpha \left(D_x \left(u * \frac{\partial g}{\partial x} \right) + D_y \left(u * \frac{\partial g}{\partial y} \right) \right)$$

$$= S + \alpha (D_x S_x + D_y S_y)$$

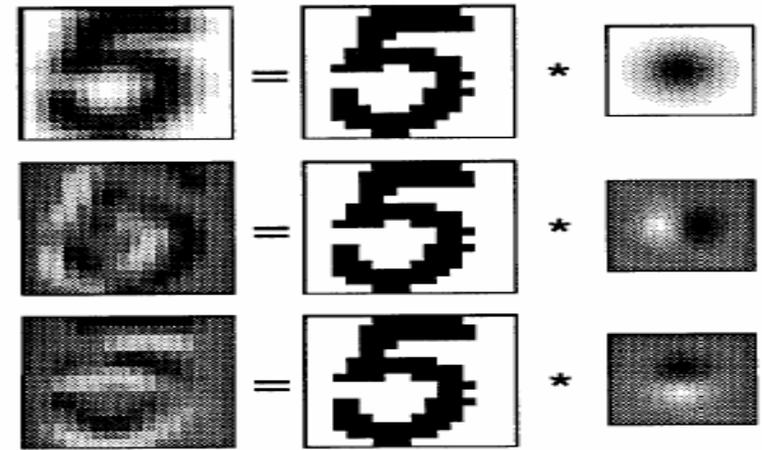


Figure 3: Graphic illustration of the computation of the three smooth images $S = u * g$, $S_x = u * \frac{\partial g}{\partial x}$ and $S_y = u * \frac{\partial g}{\partial y}$. The Gaussian function $g(x, y) = \exp^{-\frac{x^2+y^2}{2\sigma^2}}$ has a standard deviation of $\sigma = 0.9$ in this example although its graphic representation (small images on the right) have been rescaled for clarity.

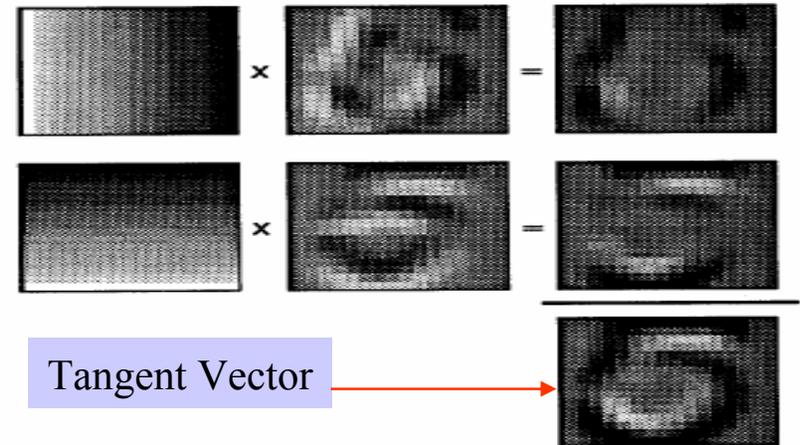
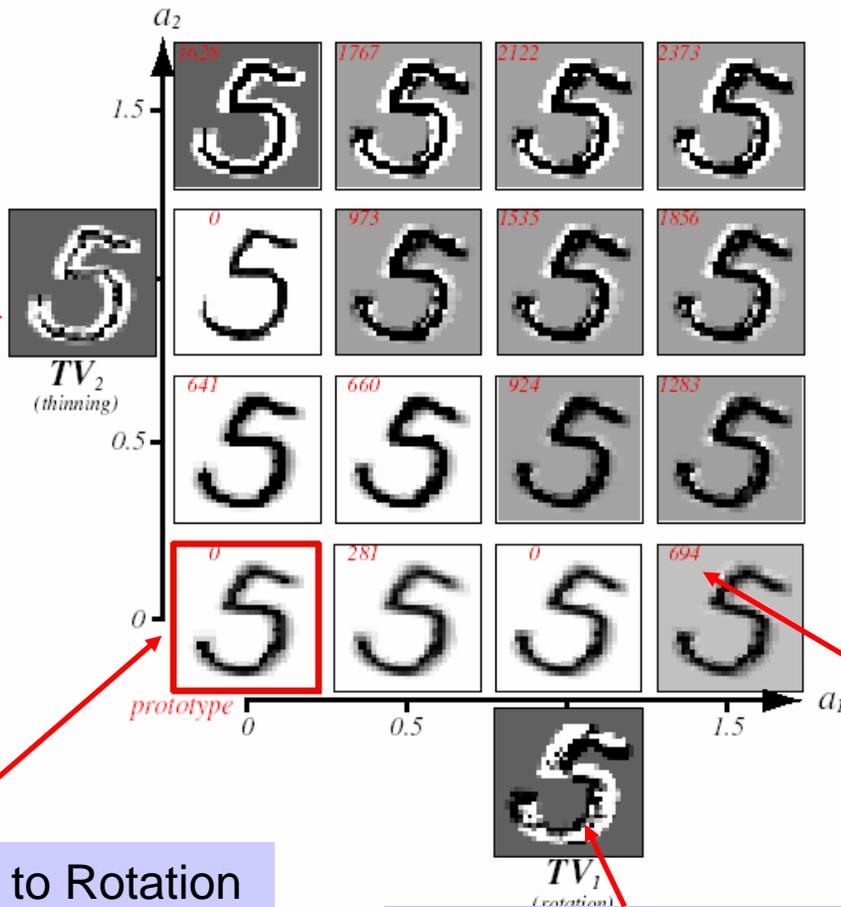


Figure 4: Graphic illustration of the computation of the tangent vector $T_u = D_x S_x + D_y S_y$ (bottom image). In this example the displacement for each pixel is proportional to the distance of the pixel to the center of the image ($D_x(x, y) = x - x_0$ and $D_y(x, y) = y - y_0$). The two multiplications (horizontal lines) as well as the addition (vertical right column) are done pixel to pixel.

Expanding a prototype using its tangent vectors

Tangent Vector
 TV_2 (thinning)



Sixteen Images obtained from Linear Combinations of prototype and two Tangent vectors TV_1 and TV_2 with coeffs a_1 and a_2

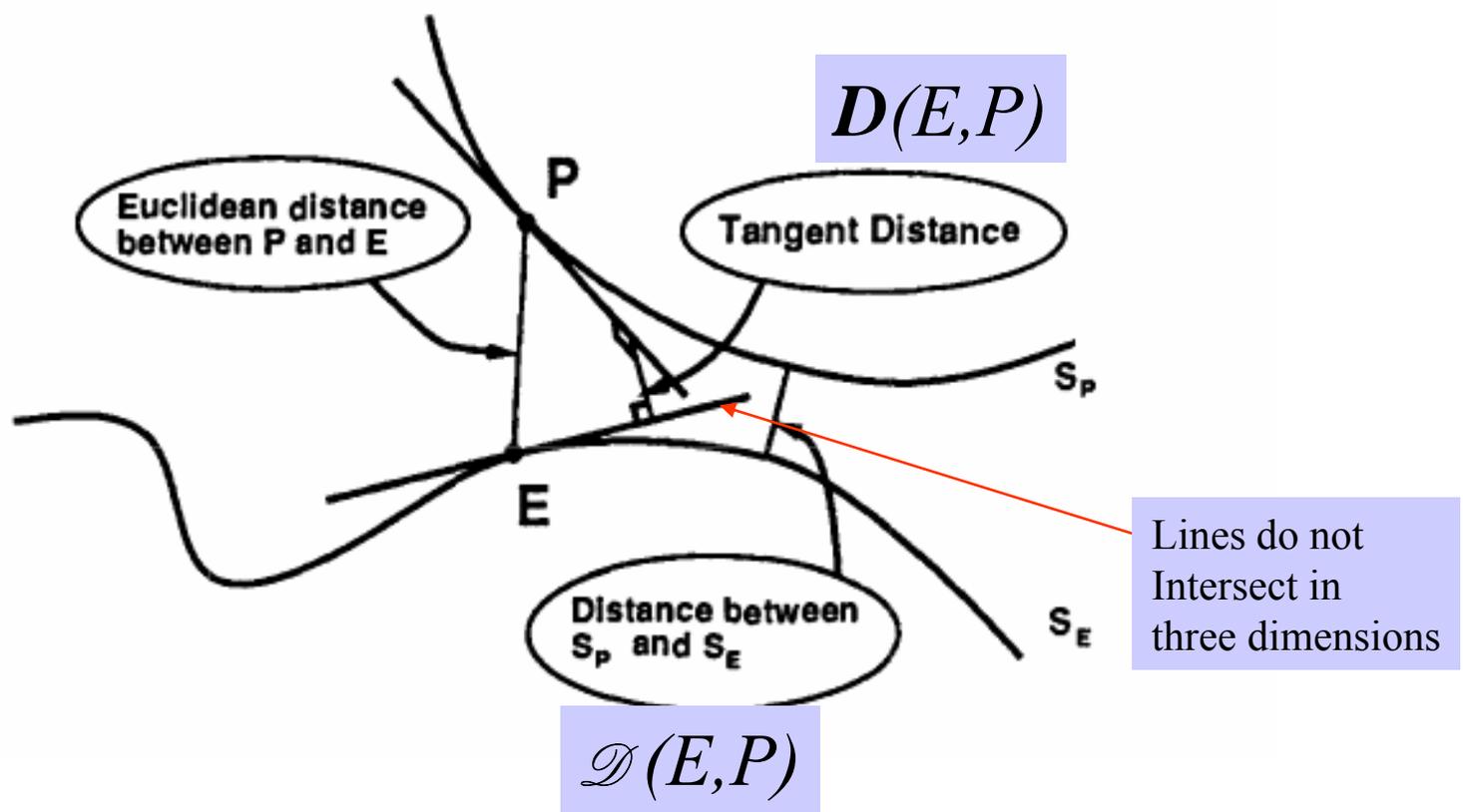
Gray value means negative

Euclidean Distance between Tangent approximation $(x' + a_1TV_1 + a_2TV_2)$ and image generated by the un-approximated transformation $F_i(x', \alpha_i)$

Prototype subjected to Rotation and Line Thinning yields two tangent vectors TV_1 and TV_2
 Computed at training time

Tangent Vector TV_1 (rotation)

Euclidean Distance, Tangent Distance, Manifold Distance



Tangent Distance: from test point x to prototype x'

- Given a matrix T consisting of the r tangent vectors at x' the tangent distance from x to x' is

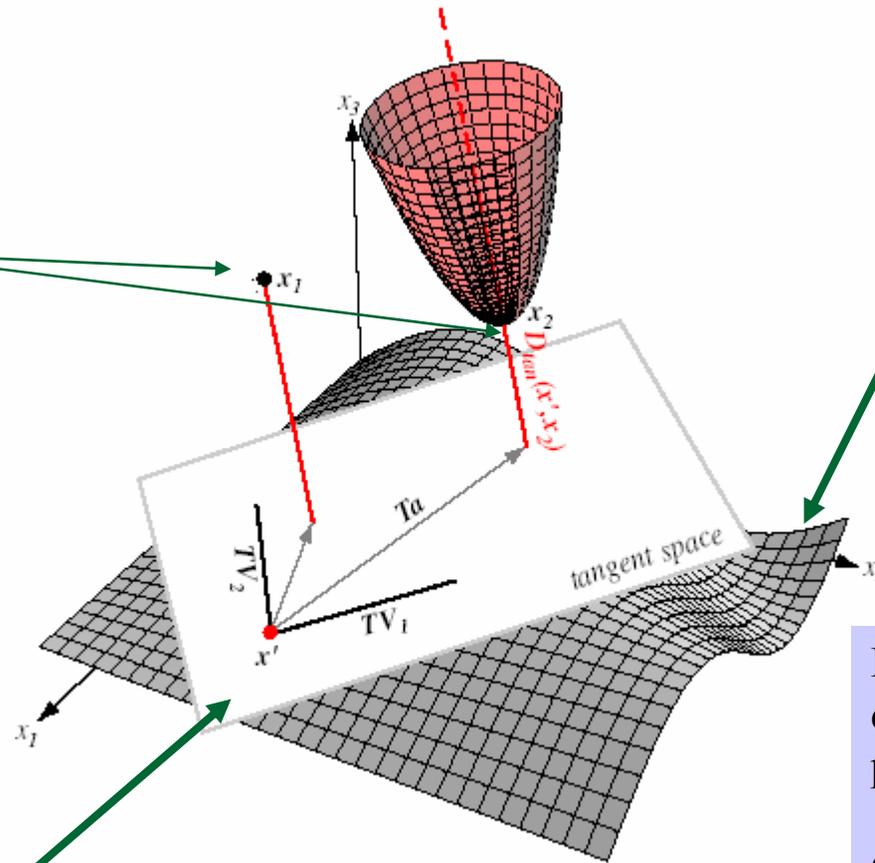
$$D_{\text{tan}}(x, x') = \min_a \left[\left\| (x' + Ta) - x \right\| \right]$$

- i.e., the Euclidean distance from x to the tangent space of x'
- This is the one-sided tangent distance since only one pattern x' is transformed
- Two-sided tangent distance improves accuracy only slightly at added computational burden
- During classification of x we find its tangent distance to x' by finding the optimizing value of a required by above equation
- Minimization is simple since squared distance we want to find is a quadratic function of a

Minimizing value a for Tangent Distance

In Euclidean space x_1 is closer to x' than x_2

In Tangent space x_2 is closer to x' than x_1



Stored prototype x' falls on this manifold when subjected to transformations

Pink paraboloid is a quadratic function of the parameter vector a

Gradient descent is used to calculate Tangent distance $D_{\tan}(x', x_2)$

Tangent Space at x' is an r -dimensional Euclidean space spanned by Tangent vectors TV_1 and TV_2