

# Tangent Distance, Tangent Prop and Manifold Tangent Classifier

Sargur N. Srihari

[srihari@cedar.buffalo.edu](mailto:srihari@cedar.buffalo.edu)

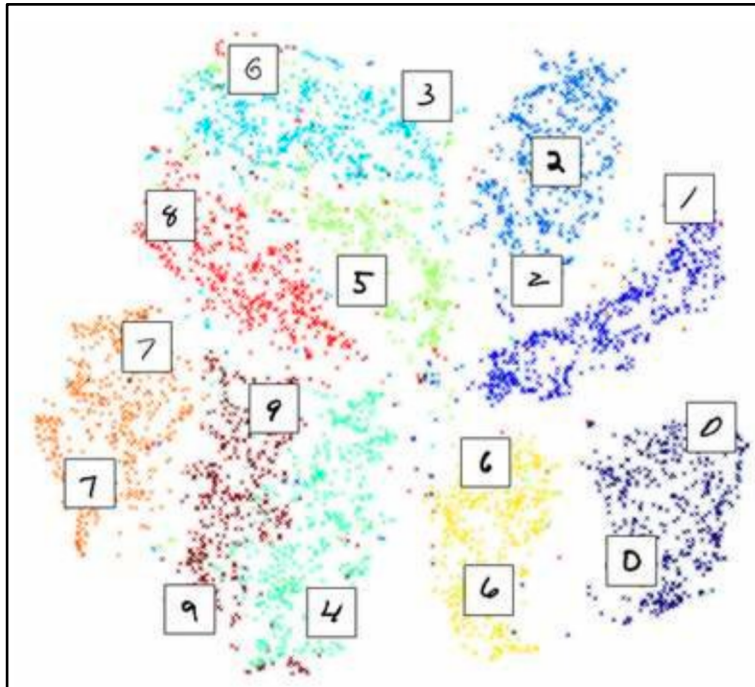
# Topics

- Manifolds in data space
- Tangent distance
- Tangent propagation

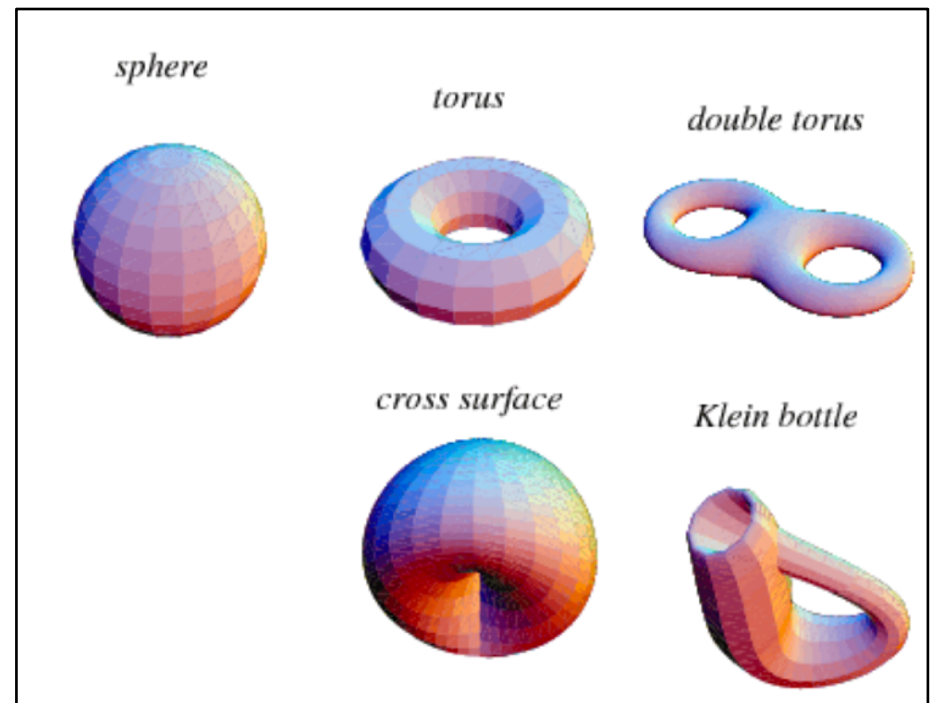
# Data lies on low dimensional manifolds

- Curse of dimensionality is overcome by the fact that data lies on low-dimensional manifolds

Manifolds of handwritten digits in 2-D



Common examples of manifolds



Note that the manifold is only the surface of these objects and not the interior

# Data manifolds are smooth

- Common theme of these examples:
  - They are somewhat smooth
    - Meaning that there are no sharp spikes or edges
- Overall shape of manifold can be amorphous
  - Describe datasets that don't have rigid boundaries



# Manifold Learning

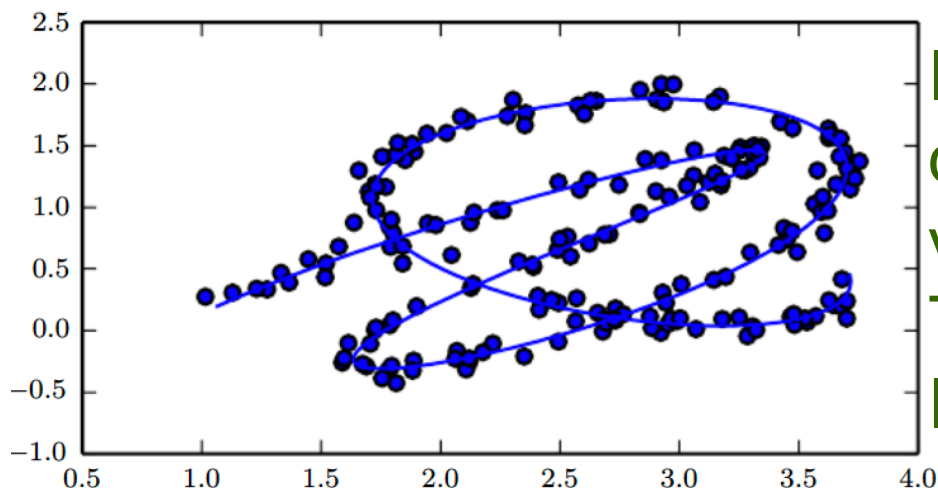
- An important idea underlying many ideas in machine learning
- A manifold is a connected region
  - Mathematically it is a set of points in a neighborhood
  - It appears to be in a Euclidean space
    - E.g., we experience the world as a 2-D plane while it is a spherical manifold in 3-D space

# Manifold in Machine Learning

- Although manifold is mathematically defined, in machine learning it is loosely defined:
  - A connected set of points that can be approximated well by considering only a small no of degrees of freedom embedded in a higher-dimensional space

Training data lying near a 1-D  
Manifold in a 2-D space

The solid line indicates the underlying  
manifold that the learner should infer



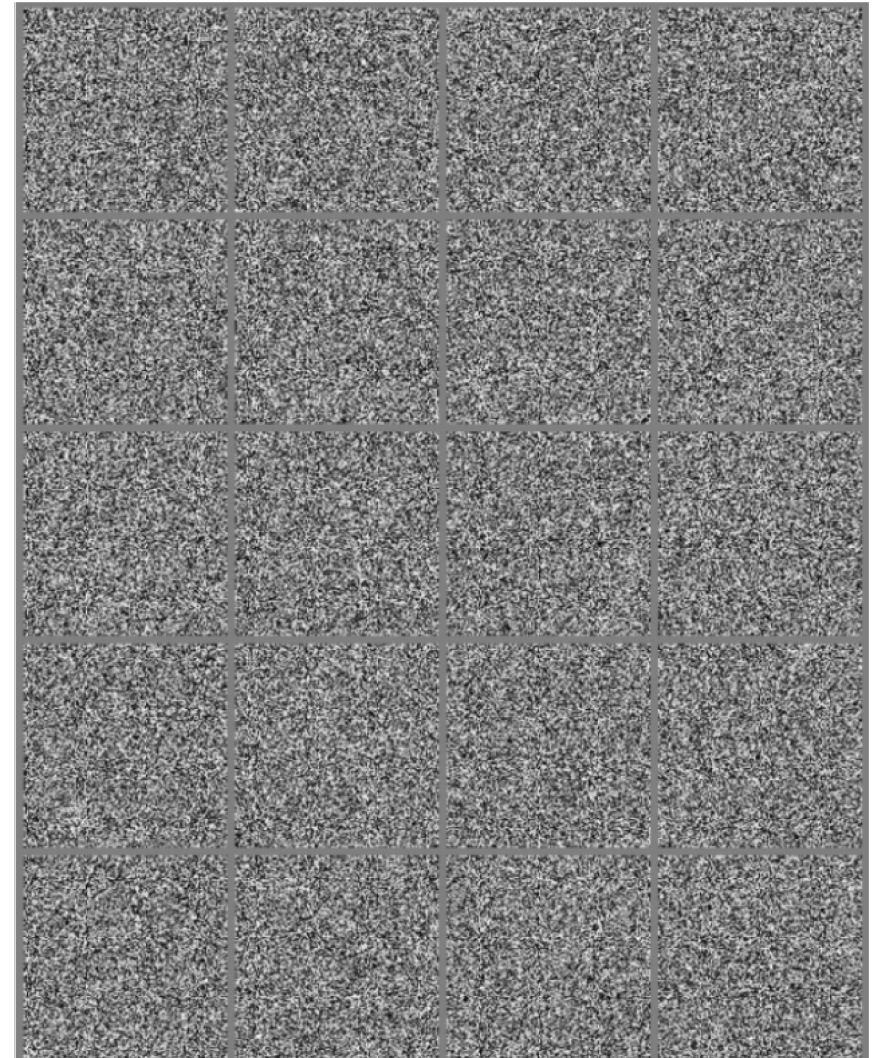
In machine learning we allow the dimensionality of the manifold to vary from one point to another. This often happens when a manifold intersects itself, as in a figure-eight

# Manifold learning surmounts $\mathbb{R}^n$

- It is sometimes hopeless to learn functions with variations across all of  $\mathbb{R}^n$
- Manifold learning algorithms surmount this obstacle by assuming most of  $\mathbb{R}^n$  consists of invalid inputs
  - And that intersecting inputs occur only along the manifolds
- Introduced for continuous data and in unsupervised learning, the probability concentration idea can be generalized to discrete and unsupervised settings

# Manifold hypothesis for Images

- Manifold assumption is justified since:
- Distributions are highly concentrated
  - Uniformly sampled points look like static noise,  
never structured
    - Although there is a non-zero probability of generating a face, it is never observed



# Manifold justified in Text domain

- If you generate a document by randomly generating text, it is a near zero probability of generating meaningful text
- Natural language sequences occupy a small volume of total space of sequences of letters

# Manifolds traced by transformations

- Manifolds can be traced by making small transformations
- Manifold structure of a dataset of human faces





# Manifolds discovered for Human Faces

- Variational autoencoder discovers underlying two-dimensional coordinate system:
  1. Rotation
  2. Emotion



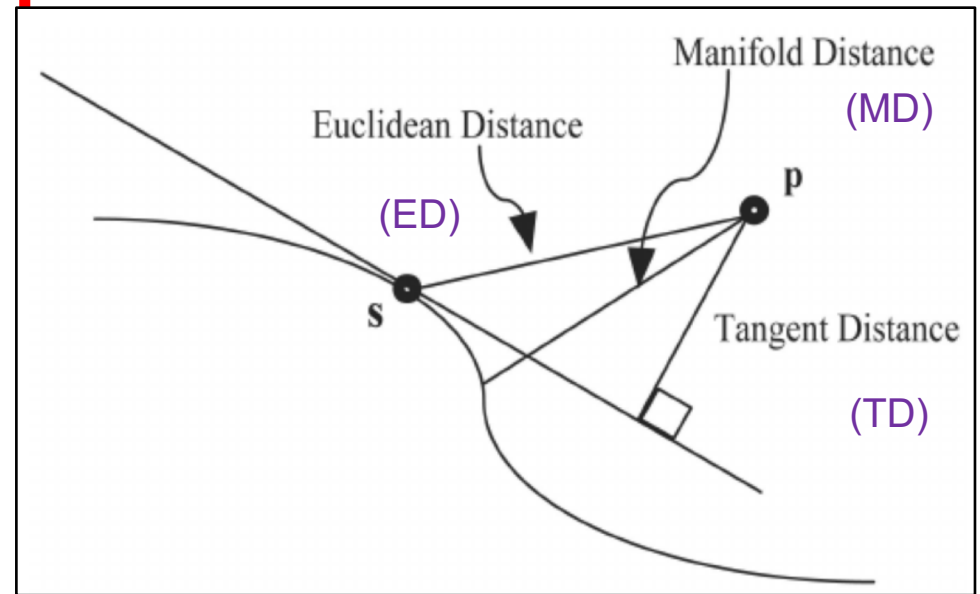
# Distance from point to manifold

Three ways to measure distance:

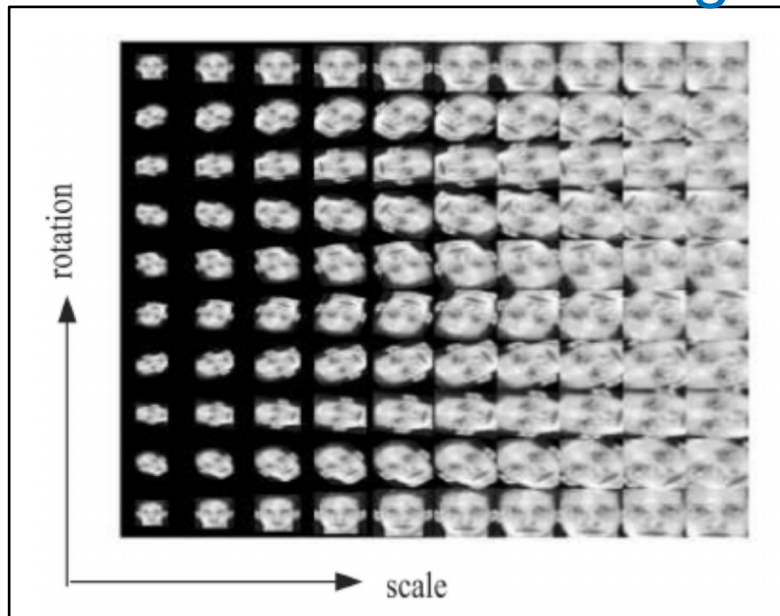
Manifold distance (MD)

Euclidean distance (ED)

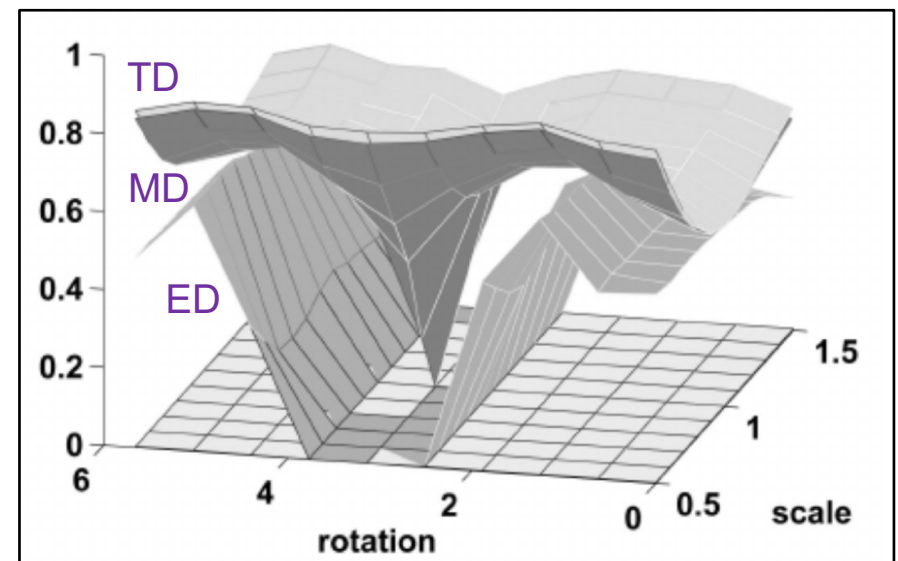
Tangent distance (TD)



Transformed face images



Distance between  $s$  and  $p$





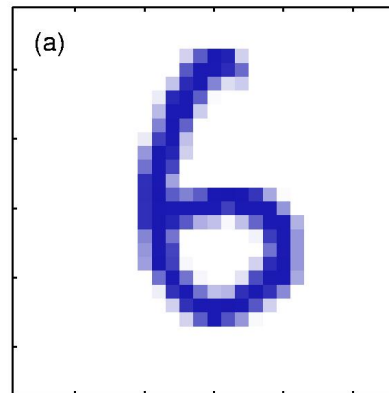
# Why tangent distance?

- Assume we are trying to classify examples
  - Also that points on a manifold have same category
- Distance between a point  $\mathbf{x}$  and a manifold is a computationally difficult task
  - Would require solving an optimization problem
- Cheap alternative is to use distance between point and tangent plane
  - By solving a linear system in the dimension of the manifold
    - Requires specifying the tangent vector

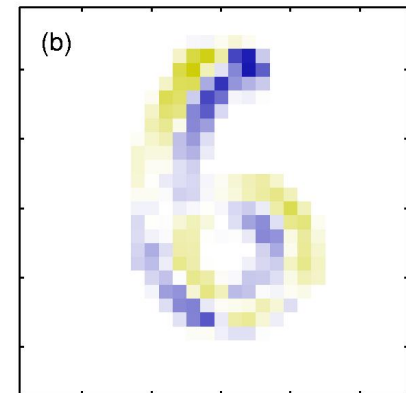
# Tangent vector from finite differences

- In practical implementation tangent vector  $\tau_n$  is approximated using finite differences by subtracting original vector  $x_n$  from the corresponding vector after transformations using a small value of  $\xi$  and then dividing by  $\xi$
- A related technique *Tangent Distance* is used to build invariance properties into distance-based methods such as *nearest-neighbor* classifiers

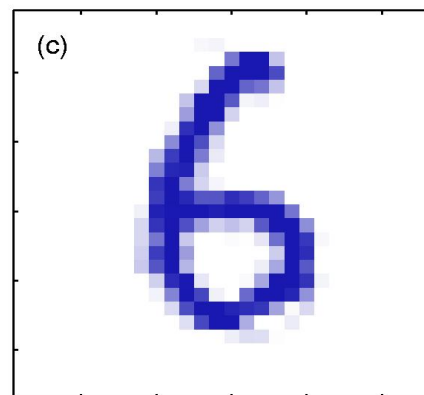
Original Image  $x$



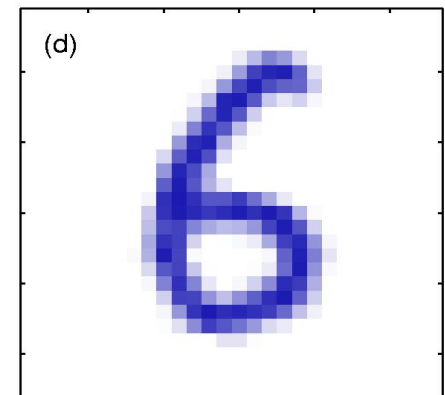
Tangent vector  $\tau$  corresponding to small clockwise rotation



Adding small contribution from tangent vector to image  $x + \epsilon\tau$

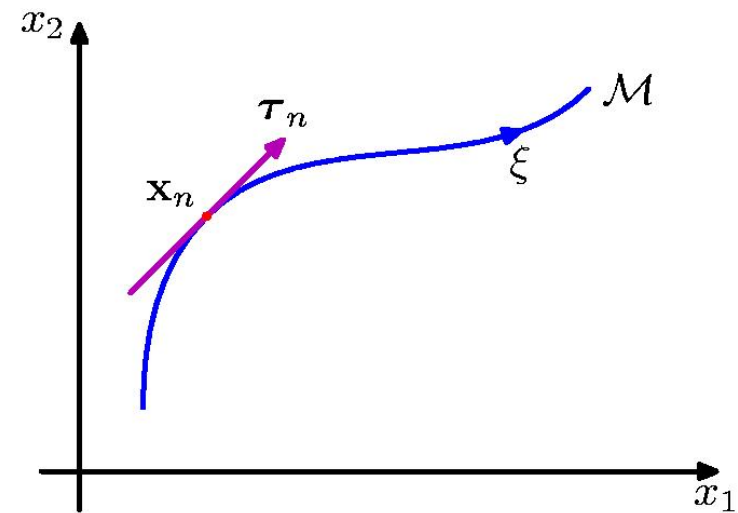


True image rotated for comparison



# Tangent Propagation

- Regularization to encourage invariance to transformations
  - by technique of tangent propagation
- Consider effect of transformation on input  $\mathbf{x}_n$ 
  - A 1- $D$  continuous transformation parameterized by  $\xi$  applied to  $\mathbf{x}_n$  sweeps a manifold  $\mathcal{M}$  in  $D$ -dimensional input space



Two-dimensional input space showing effect of continuous transformation with single parameter  $\xi$

Let the vector resulting from acting on  $\mathbf{x}_n$  by this transformation be denoted by  $s(\mathbf{x}_n, \xi)$  defined so that  $s(\mathbf{x}, 0) = \mathbf{x}$ .

Then the tangent to the curve  $\mathcal{M}$  is given by the directional derivative  $\tau = \delta s / \delta \xi$  and the tangent vector at point  $\mathbf{x}_n$  is given by

$$\tau_n = \left. \frac{\partial s(\mathbf{x}_n, \xi)}{\partial \xi} \right|_{\xi=0}$$

# Tangent Prop as Regularization

- Under a transformation of input vector
  - The network output vector will change
  - Derivative of output  $k$  wrt  $\xi$  is given by

$$\left. \frac{\partial y_k}{\partial \xi} \right|_{\xi=0} = \sum_{i=1}^D \frac{\partial y_k}{\partial x_i} \left. \frac{\partial x_i}{\partial \xi} \right|_{\xi=0} = \sum_{i=1}^D J_{ki} \tau_i$$

- where  $J_{ki}$  is the  $(k, i)$  element of Jacobian Matrix  $J$
- Result used to modify standard error function

$$\begin{aligned} \tilde{E} &= E + \lambda \Omega \\ \text{where } \lambda &\text{ is a regularization coefficient and} \\ \Omega &= \frac{1}{2} \sum_n \sum_k \left( \left. \frac{\partial y_{nk}}{\partial \xi} \right|_{\xi=0} \right)^2 = \frac{1}{2} \sum_n \sum_k \left( \sum_{i=1}^D J_{nki} \tau_{ni} \right)^2 \end{aligned}$$

- To encourage local invariance in neighborhood of data pt