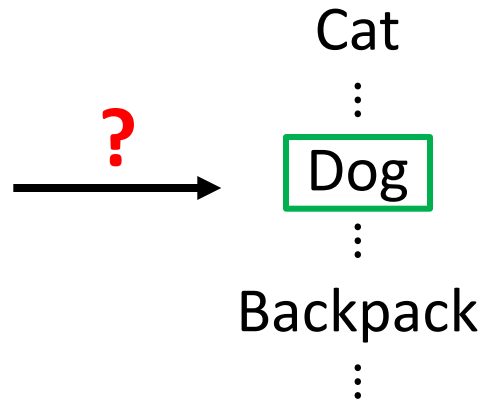


Differential Geometric Regularization for Supervised Learning of Classifiers

Qinxun Bai
Boston University

Visual Recognition



Supervised learning of classifiers

- State-of-the-art on ImageNet Challenge

Human level: classification error $< 4\%$

Counter-Intuitive Properties

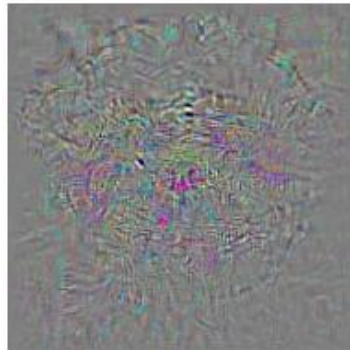
Dog?

Yes



training

+



=



testing

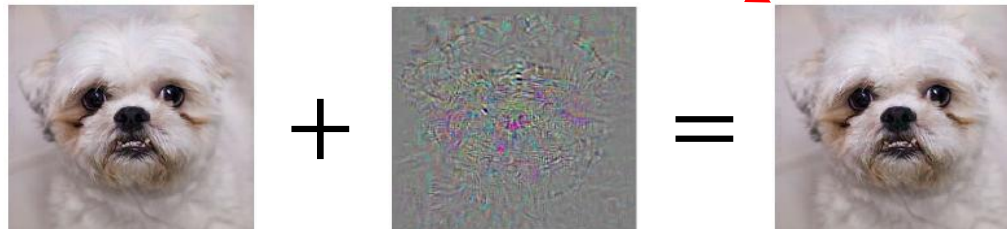
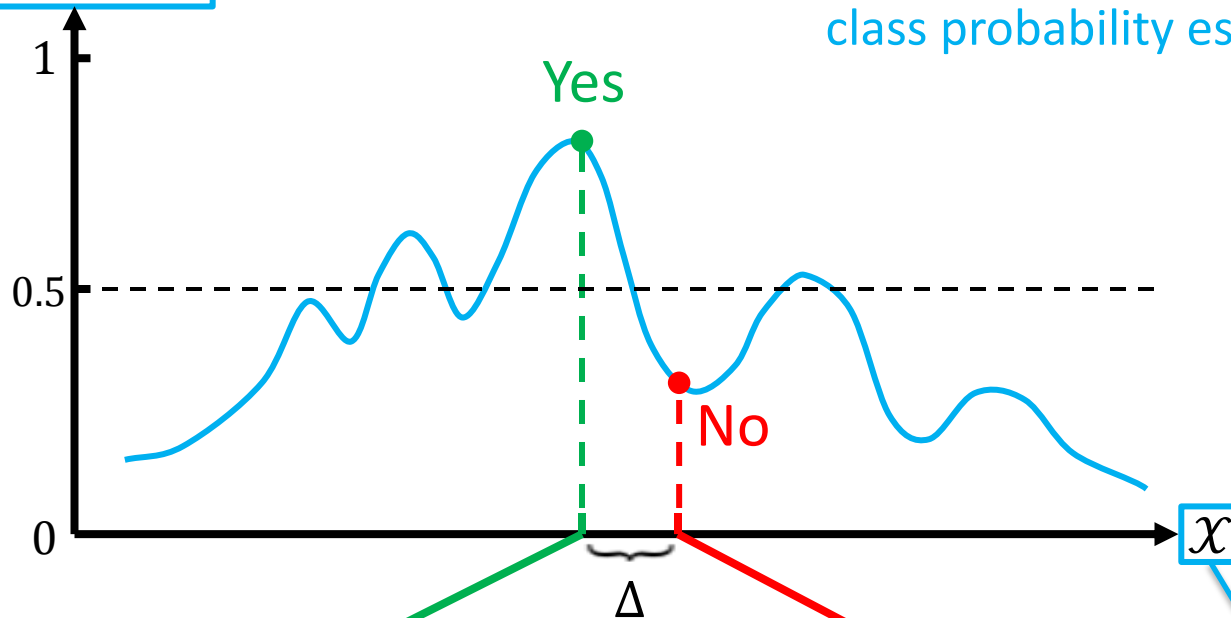
No

Fool DNN by hardly perceptible perturbation [Szegedy *et al.* 2013]

Rapid Local Oscillation

class probability

$$P(y = 1 | \mathbf{x})$$

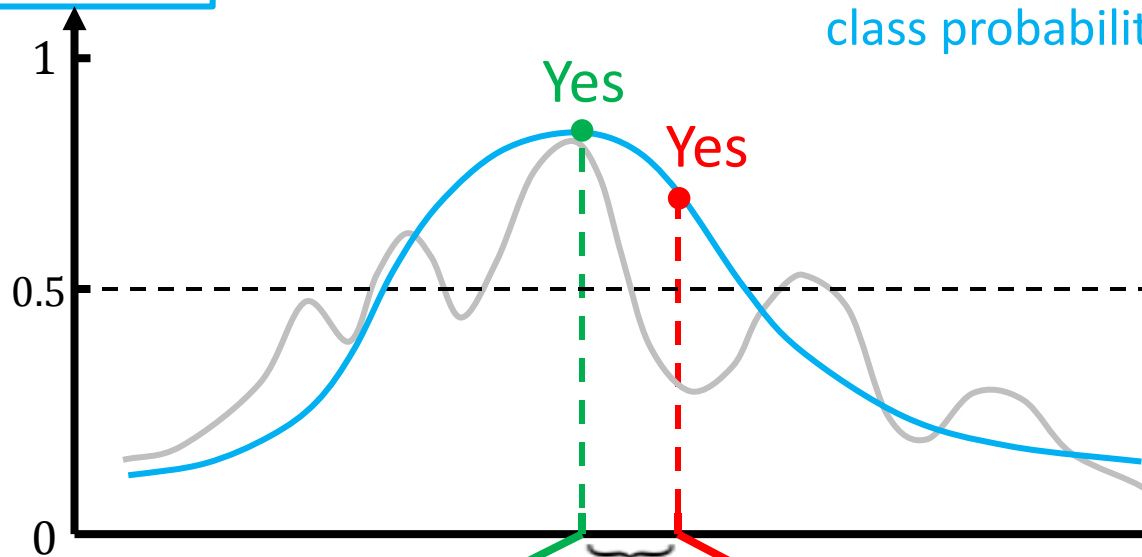


space of images
(high-dimensional)

Rapid Local Oscillation

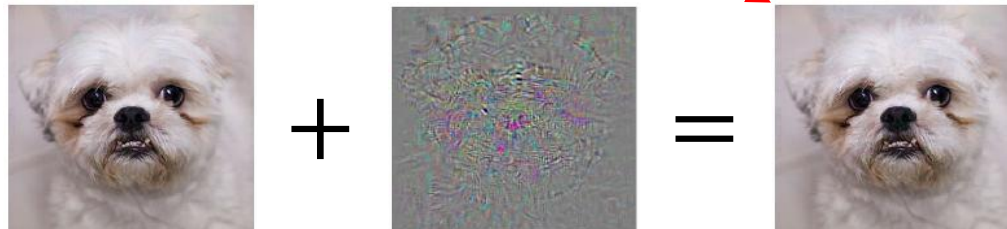
class probability

$$P(y = 1 | \mathbf{x})$$

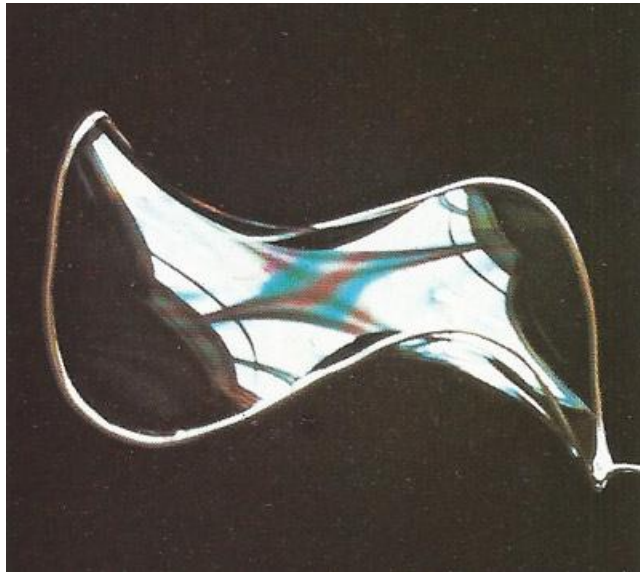


\mathbf{x}

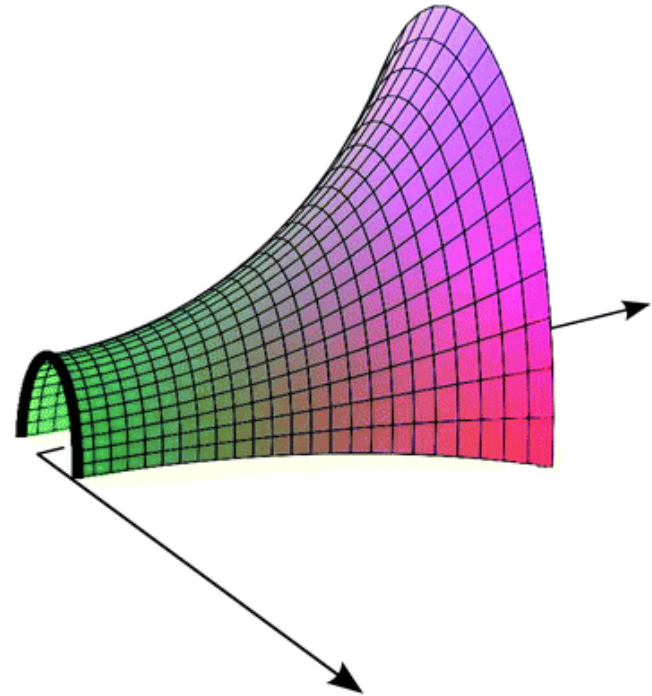
space of images
(high-dimensional)



Geometric Idea: Minimal Surfaces



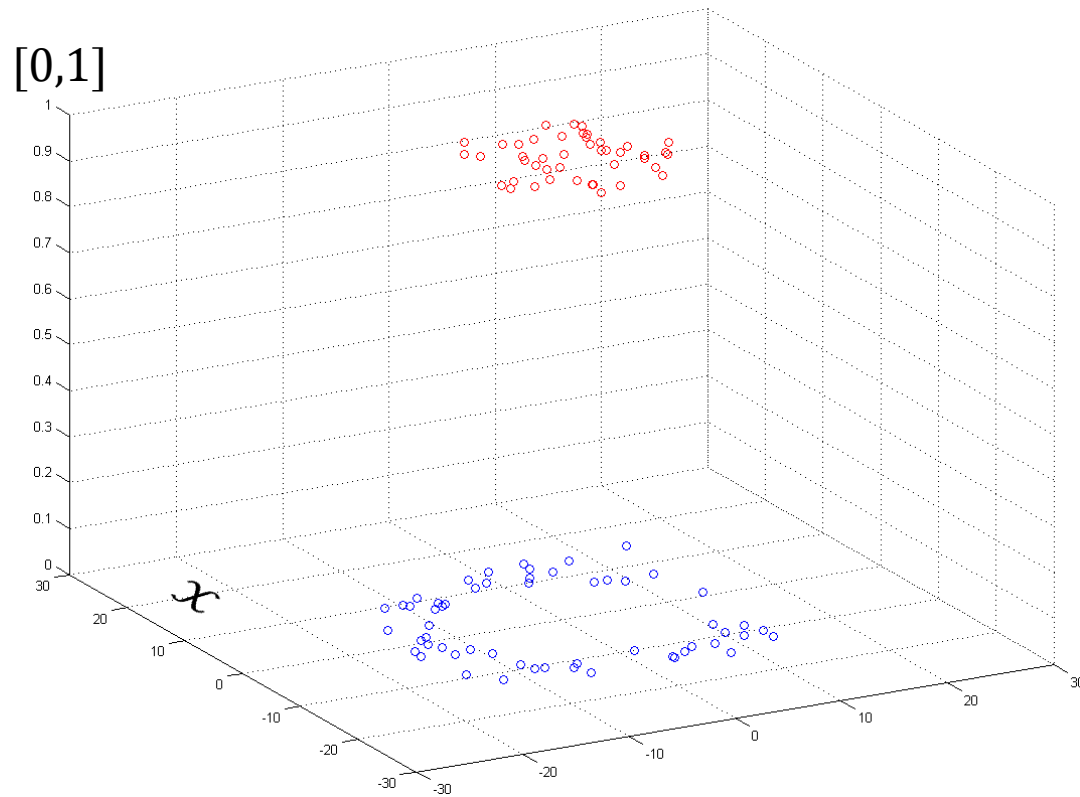
soap film



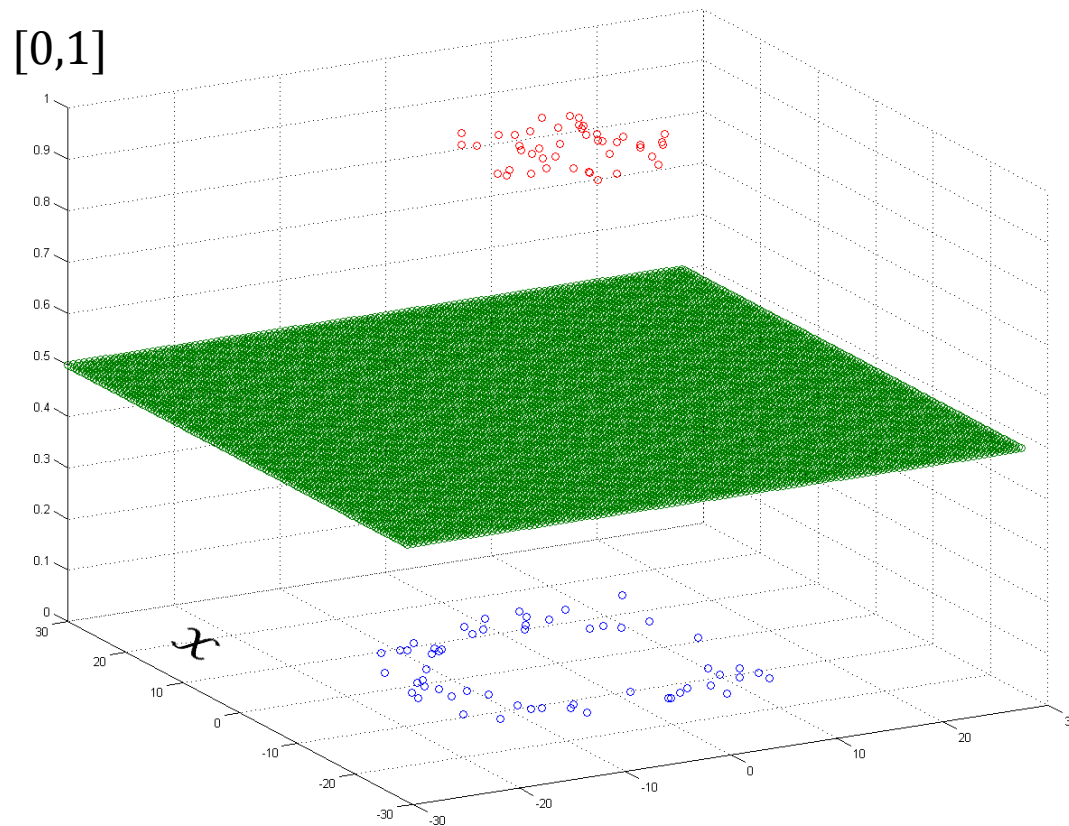
catenoid

image credit: Google image search

Physical Model

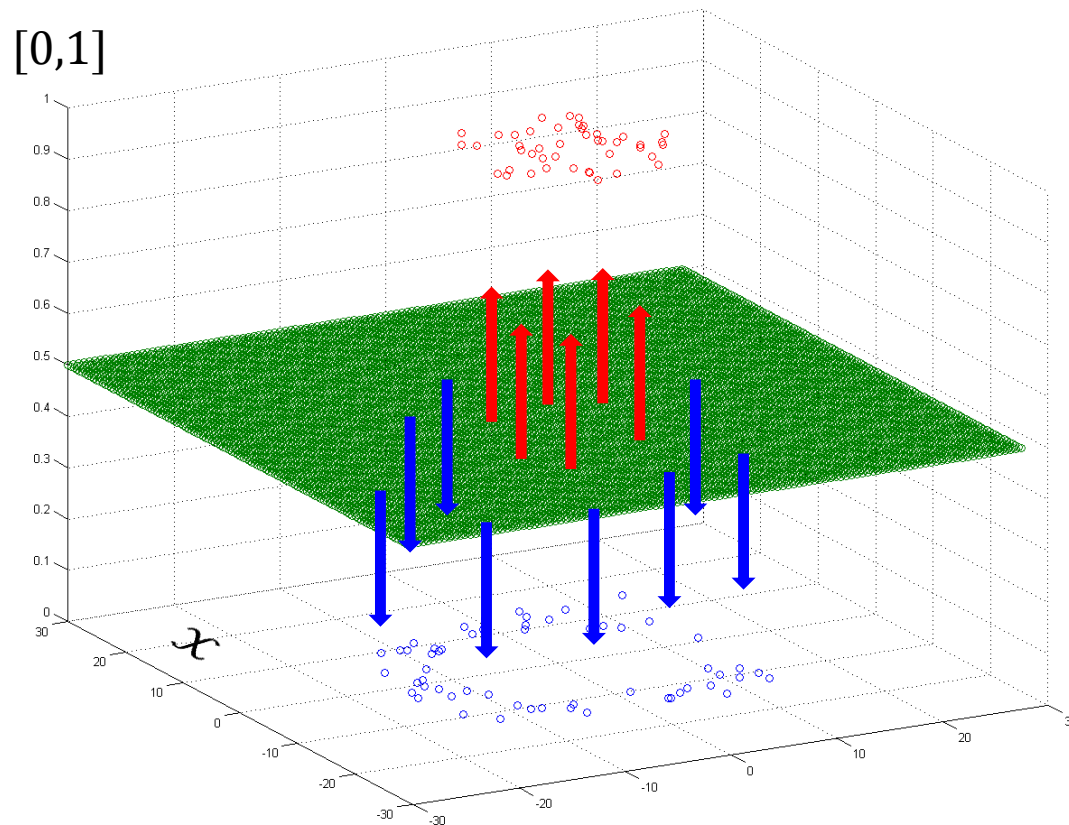


Physical Model



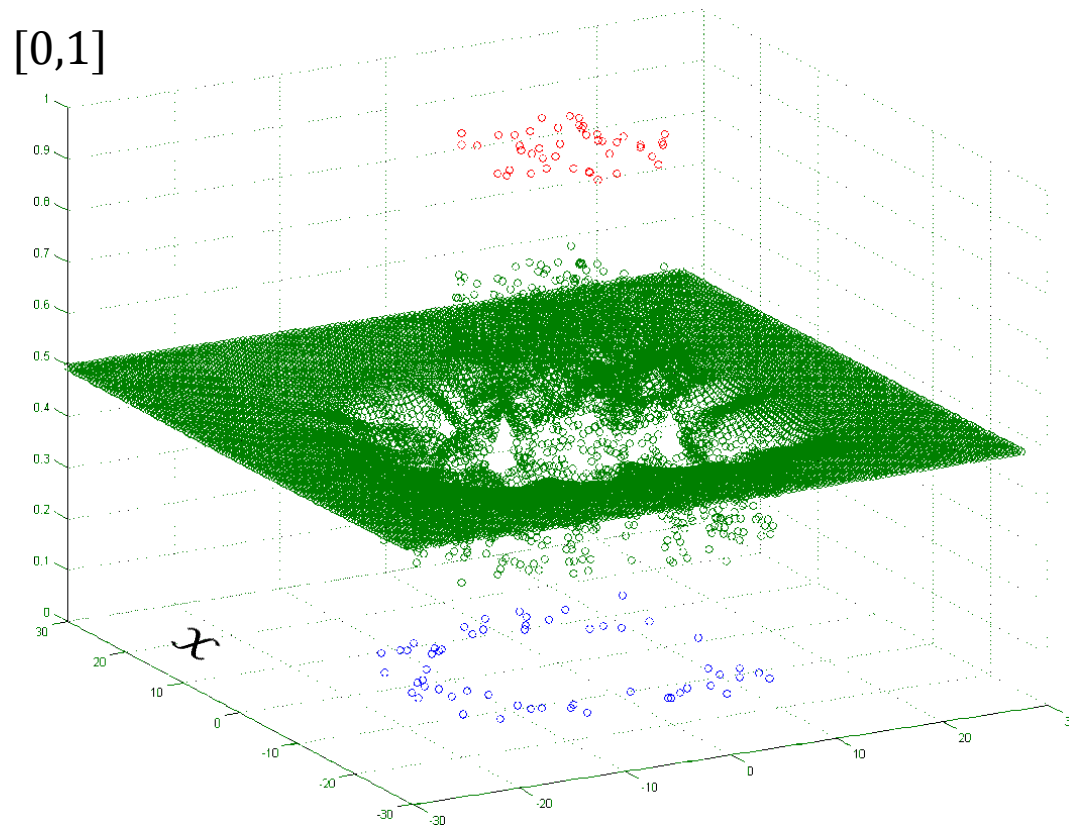
Initial hyper-surface

Physical Model



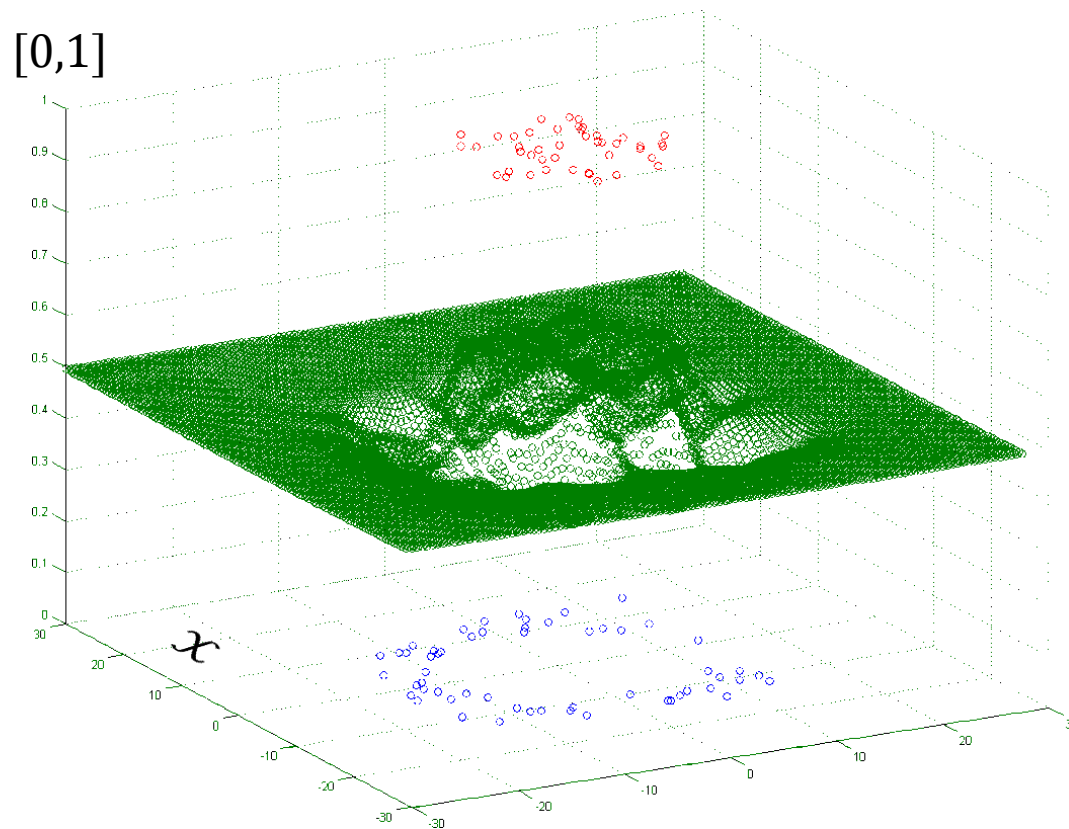
hyper-surface deforms towards training data
as if attracted by gravitational force due to point masses

Physical Model



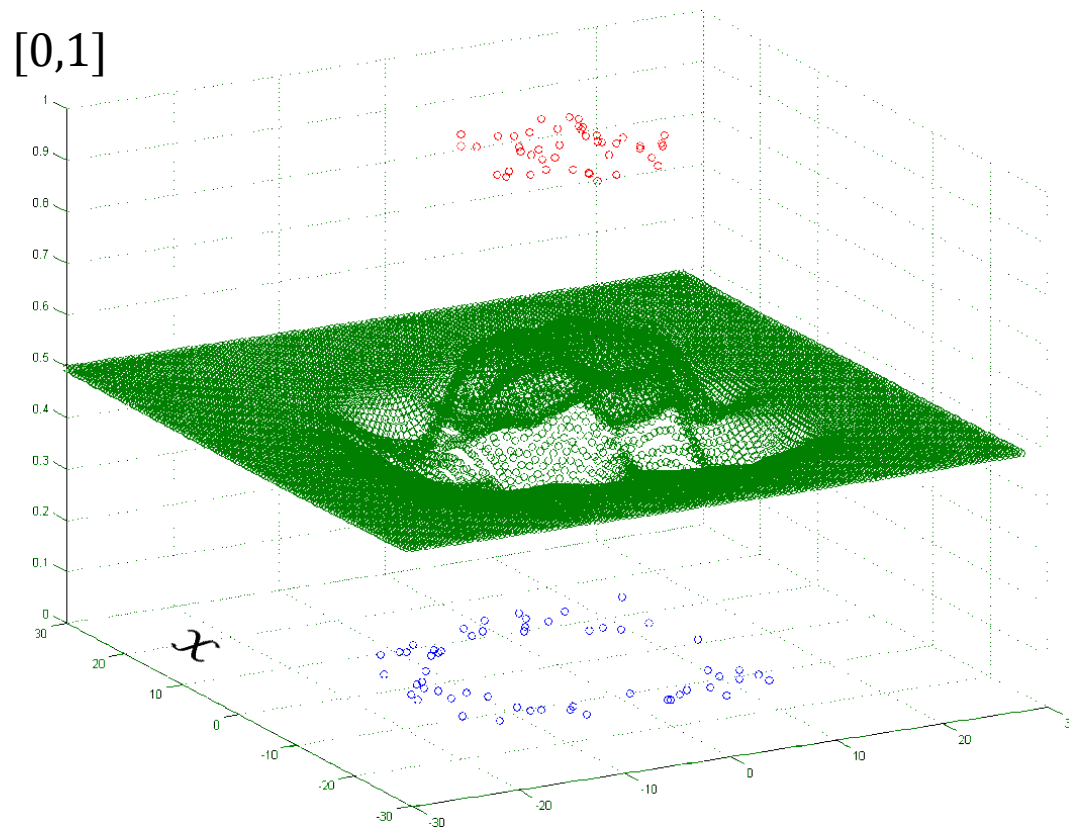
hyper-surface deforms towards training data
as if attracted by gravitational force due to point masses

Physical Model



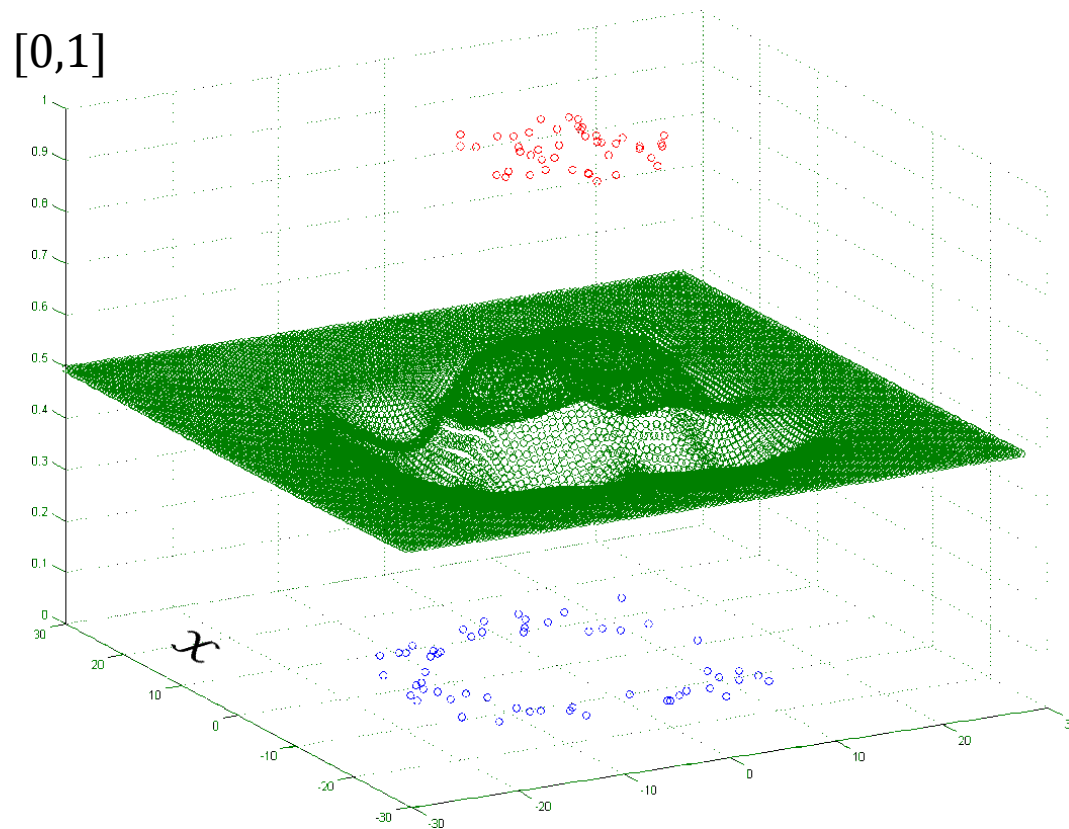
hyper-surface remains as tight as possible
as if in the presence of surface tension

Physical Model



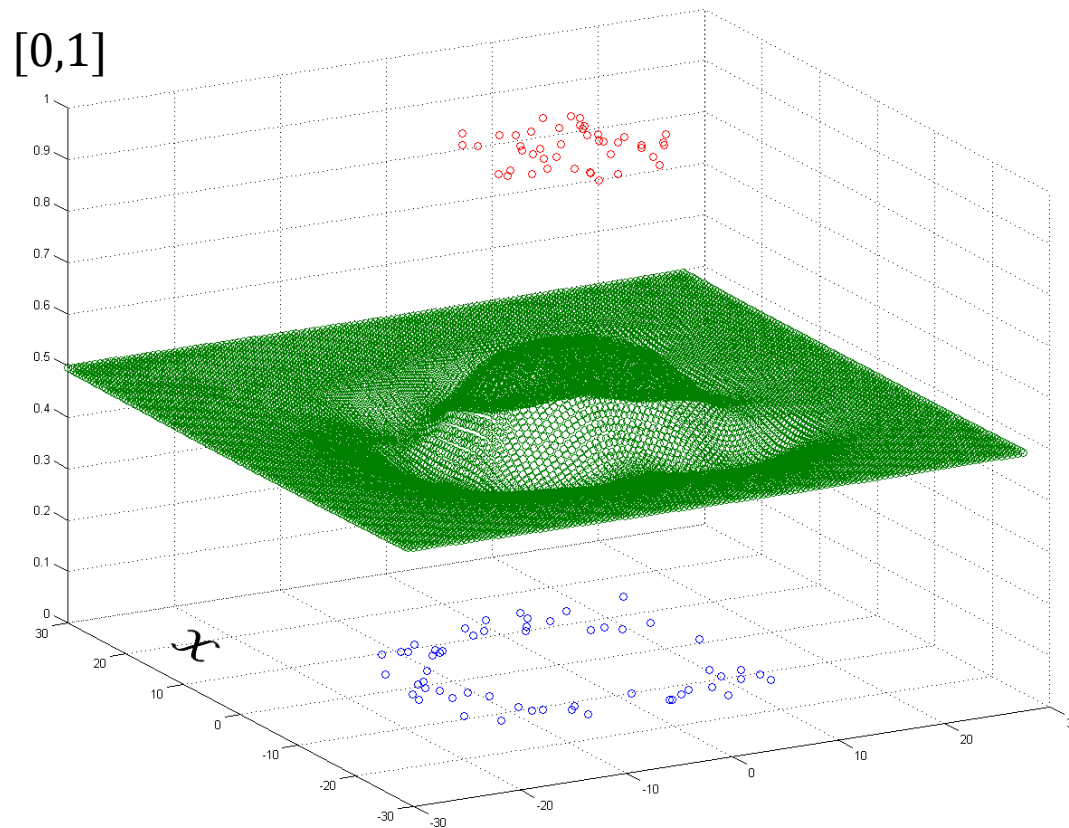
hyper-surface remains as tight as possible
as if in the presence of surface tension

Physical Model



hyper-surface remains as tight as possible
as if in the presence of surface tension

Physical Model



hyper-surface remains as tight as possible
as if in the presence of surface tension

Formal Setup

Learn a function $f: \mathcal{X} \rightarrow \Delta^K$ as an estimator of $P(y|\mathbf{x})$

input feature
space

output probabilistic
simplex for K classes

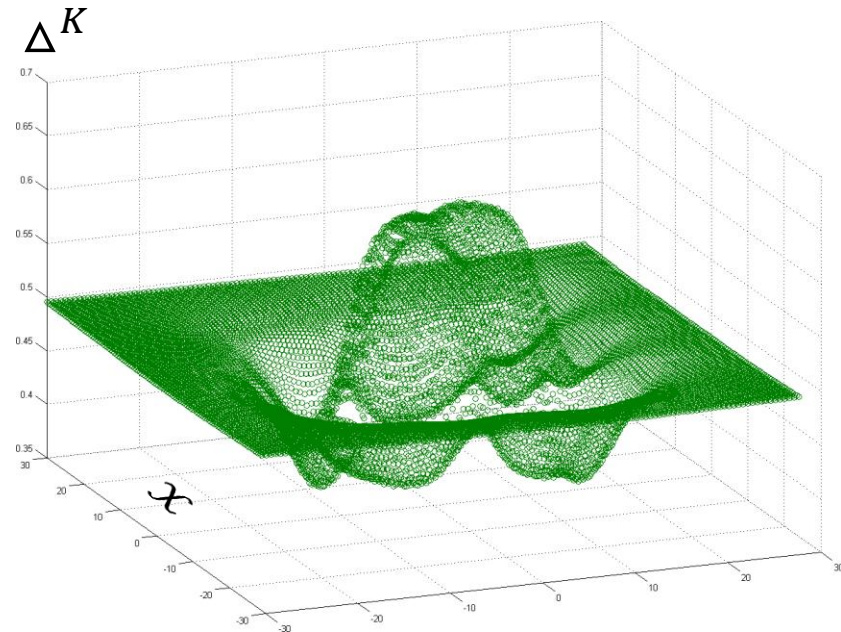
Formal Setup

Learn a function $f: \mathcal{X} \rightarrow \Delta^K$ as an estimator of $P(y|\mathbf{x})$

Hyper-surface associated with f :

$$\mathit{graph}(f) = \{(\mathbf{x}, f^1(\mathbf{x}), \dots, f^K(\mathbf{x})) | \mathbf{x} \in \mathcal{X}\} \in \mathcal{X} \times \Delta^K$$

exploit the geometry
of this hyper-surface!



Regularization Scheme

Minimize the regularized loss \mathcal{P} in functional space \mathcal{H}

$$\min_{f \in \mathcal{H}} \mathcal{P}(f) = \min_{f \in \mathcal{H}} \{L(f) + \lambda G(f)\}$$

Data term

penalize the error of f in explaining the training data

Regularization term

penalize the volume of $graph(f)$

Effect of Regularization

Optimization perspective

$$\min_{f \in \mathcal{H}} (L(f) + \lambda G(f)) \Leftrightarrow \min_{f \in \mathcal{H}_\lambda} L(f)$$

functional space that
the algorithm works in

shrunk functional space
 $\mathcal{H}_\lambda = \{f \in \mathcal{H}, G(f) \text{ bounded}\}$

- imposing $G(f) \Leftrightarrow$ shrinking $\mathcal{H} \rightarrow \mathcal{H}_\lambda$
- properly-shrink is the key for generalization
- sculpturing: λ is your hand, $G(f)$ is the knife!

Shrink the Search Space \mathcal{H}

Decomposition of excess error:

$$\underbrace{R(f)}_{\text{generalization risk}} - \underbrace{R(f^*)}_{\text{Bayes risk (optimal)}} = \underbrace{(R(f) - R(\mathcal{H}))}_{\text{optimal risk achievable in } \mathcal{H}} + (R(\mathcal{H}) - R(f^*))$$

Shrink the Search Space \mathcal{H}

Decomposition of excess error:

$$R(f) - R(f^*) = \underbrace{(R(f) - R(\mathcal{H}))}_{\text{estimation error}} + \underbrace{(R(\mathcal{H}) - R(f^*))}_{\text{approximation error}}$$

Concerning estimation error:

smaller \mathcal{H} , smaller estimation error

Concerning approximation error:

larger \mathcal{H} , smaller approximation error

A subtle trade-off: proper shrinking of \mathcal{H} , ideally:

$G(f)$ precisely encodes our prior & cross-validation on λ

Functional-norms: Smoothness

Functional perspective

- penalizing functional norm \rightarrow smoothness

Smoothness of different kinds:

- not specifically tailored to measure the amount of local oscillation
- overkill the hypothesis space
- Sculpturing with an axe? Need a sculptor's knife!

Our Argument: Mean Curvature

Geometric perspective:

- $\{(\mathbf{x}, f(\mathbf{x})) | \mathbf{x} \in \mathcal{X}\}$: a submanifold in $\mathcal{X} \times \Delta^K$

Mean Curvature of this submanifold:

- in differential geometric sense
- a specific measure of the amount of local oscillation
- generalizes to high dimensional space
- handles binary and multiclass uniformly

Existing Geometric Regularization

On geometry of the **marginal distribution** $P(\mathbf{x})$

- manifold regularization (Belkin *et al.* 2006)

On geometry of the **decision boundary** in \mathcal{X}

- level set based regularization

Cai & Sowmya 2007; Varshney & Willsky 2010

- Euler's Elastica based regularization

Lin *et al.* 2012; 2015

The small local oscillation of $\eta(\mathbf{x})$ is **Not captured**

Solve for $\min_{f \in \mathcal{H}} \mathcal{P}(f)$

Solving it directly is too difficult!

Solve iteratively by gradient flow: $\frac{df_t}{dt} = -\nabla \mathcal{P}$

- starting from neutral estimator $f_0 = \left(\frac{1}{K}, \dots, \frac{1}{K}\right)$
- evolve f_t towards $-\nabla \mathcal{P}$
- f_t will flow to a local minimum of \mathcal{P}

Algorithm (binary & multiclass)

Input: training data, trade-off λ , step-size τ

Initialize: $f(\mathbf{x}_i; \mathbf{w}) = \left(\frac{1}{K}, \dots, \frac{1}{K}\right)$, $M = \frac{\partial f}{\partial \mathbf{w}}$

For $t = 1$ to T

- Evaluate gradient vector $\nabla \mathcal{P}$ at every training point \mathbf{x}_i
- $\mathbf{w} \leftarrow \mathbf{w} - \tau M^{-1} [\nabla \mathcal{P}(\mathbf{x}_1), \dots]^T$

Output: class probability f

Solid math
Simple algorithm
Parallelizable!

Geometric Foundation on \mathcal{H}

$$\mathcal{H} = \text{Maps}(\mathcal{X}, \Delta^K), \mathcal{H}' = \text{Maps}(\mathcal{X}, \mathbb{R}^K)$$

Topology

- Frechet topology on \mathcal{H}' , and the induced topology on \mathcal{H}
i.e. two functions in \mathcal{H} are close if the functions and all their partial derivatives are pointwise close

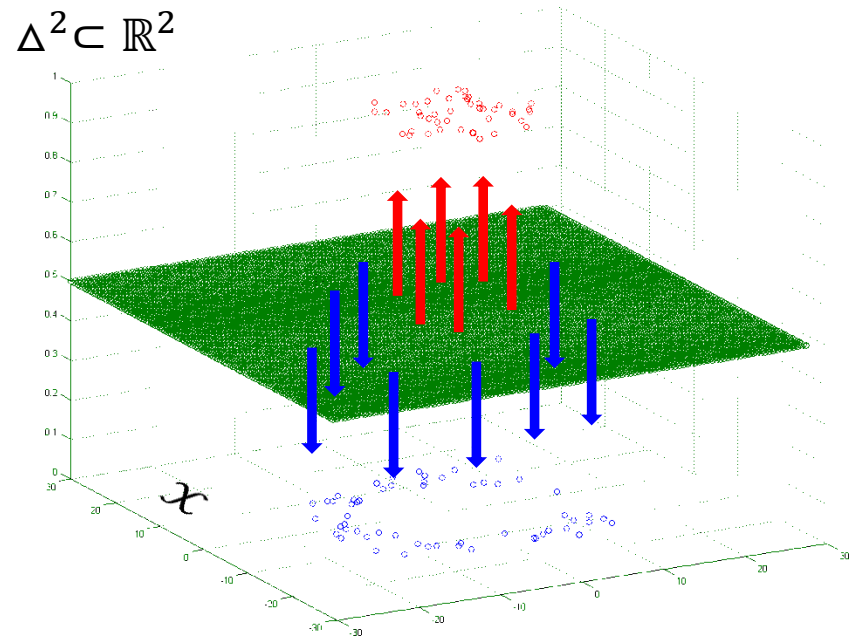
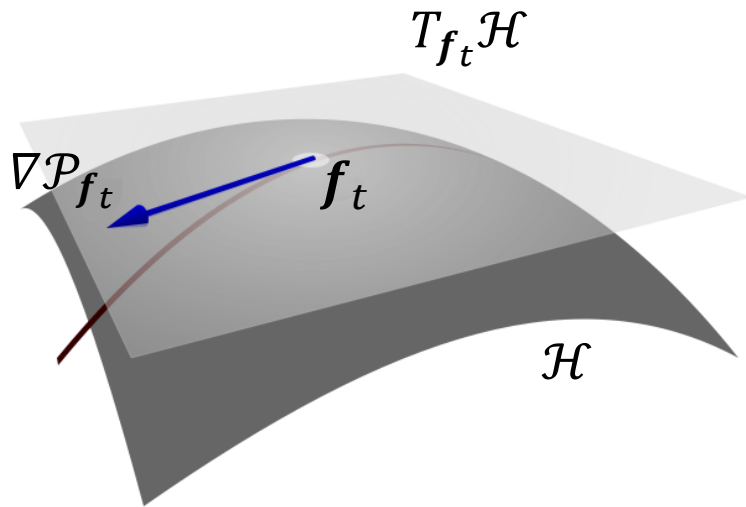
Riemannian metric

- Restrict the L^2 metric on \mathcal{H}' to each tangent space $T_f\mathcal{H}$

$$\langle \phi_1, \phi_2 \rangle = \int_{\mathcal{X}} \phi_1(\mathbf{x})\phi_2(\mathbf{x})dvol_{\mathbf{x}}$$

where $\phi_i \in \mathcal{H}'$ and $dvol_{\mathbf{x}}$ is the volume form of the induced Riemannian metric on $graph(\mathbf{f})$.

The Gradient $\nabla \mathcal{P}_{f_t}$



$\nabla \mathcal{P}_{f_t}$: tangent vector in $T_{f_t} \mathcal{H}$ \longleftrightarrow vector field on $graph(f_t)$

Computation of $\nabla \mathcal{P} = \nabla L + \lambda \nabla G$

Computing ∇L is easy

- *e.g.* back propagation for neural networks

Computing ∇G : mean curvature flow

- **Our Theorem:**

need **only 1st and 2nd** partial derivatives of f , rest of computation is just matrix manipulations

Empirical Data Term

Quadratic loss

$$L(\mathbf{f}) = \sum_{i=1}^m \|f(\mathbf{x}_i) - \mathbf{z}_i\|^2$$

Cross-entropy loss

$$L(\mathbf{f}) = - \sum_{i=1}^m \sum_{l=1}^K z_i^l \log f^l(\mathbf{x}_i)$$

computation of $\nabla L(\mathbf{x}_i)$ is trivial for both losses

Geometric Regularization Term

Volume penalty

$$G(\mathbf{f}) = \int_{\text{graph}(\mathbf{f})} d\text{vol} = \int_{\text{graph}(\mathbf{f})} \sqrt{\det(\mathbf{g})} dx^1 \cdots dx^N$$

\mathbf{g} is the Riemannian metric on $\text{graph}(\mathbf{f})$ induced from the standard dot product on \mathbb{R}^{N+K}

Geometric Regularization Term

Gradient vector field of $G(\mathbf{f})$

$$-\nabla G = \text{Tr} \mathbb{H}^K$$

$$= (g^{-1})^{ij} (f_{ji}^1 - (g^{-1})^{rs} f_{rs}^l f_i^l f_j^1, \dots, f_{ji}^K - (g^{-1})^{rs} f_{rs}^l f_i^l f_j^K)$$

where f_i^l, f_{ij}^l denote partial derivatives of f^l

given 1st and 2nd partial derivatives

the computation involves only matrix manipulations

Example Formulation: RBFs

Represent f as “softmax” output of RBFs

$$f^j = \frac{\exp(h^j)}{\sum_{l=1}^K \exp(h^l)}, \quad h^j = \sum_{i=1}^m a_i^j \varphi_i(\mathbf{x}), \quad \text{for } j = 1, \dots, K$$

where $\varphi_i(\mathbf{x}) = e^{-\frac{1}{c}\|\mathbf{x}-\mathbf{x}_i\|^2}$ is the RBF centered at \mathbf{x}_i

Gradient update for $A = (a_i^l)$

$$A \leftarrow A - \tau M^{-1} [\nabla \mathcal{P}_h(\mathbf{x}_1), \dots, \nabla \mathcal{P}_h(\mathbf{x}_m)]^T,$$

where $\nabla \mathcal{P}_h(\mathbf{x}_i) = \left[\frac{\partial f}{\partial \mathbf{h}} \right]_{\mathbf{x}_i}^T \nabla \mathcal{P}_f(\mathbf{x}_i)$, $M_{ij} = \varphi_j(\mathbf{x}_i)$

Experiments – RBF Representation

Datasets from UCI Repository

- Four binary and four multiclass datasets
- Following the choice/setup of previous papers

Comparing with two groups of classifiers

- RBF + functional norm regularization: RBN, SVM, KLR
- RBF + existing geometric regularization: LLS, GLS, EE

UCI Datasets – Interesting pairs

KLR vs. Ours-CE

- same: RBF-based, cross-entropy loss
- diff regularizer: RKHS norm vs geometry on class probability

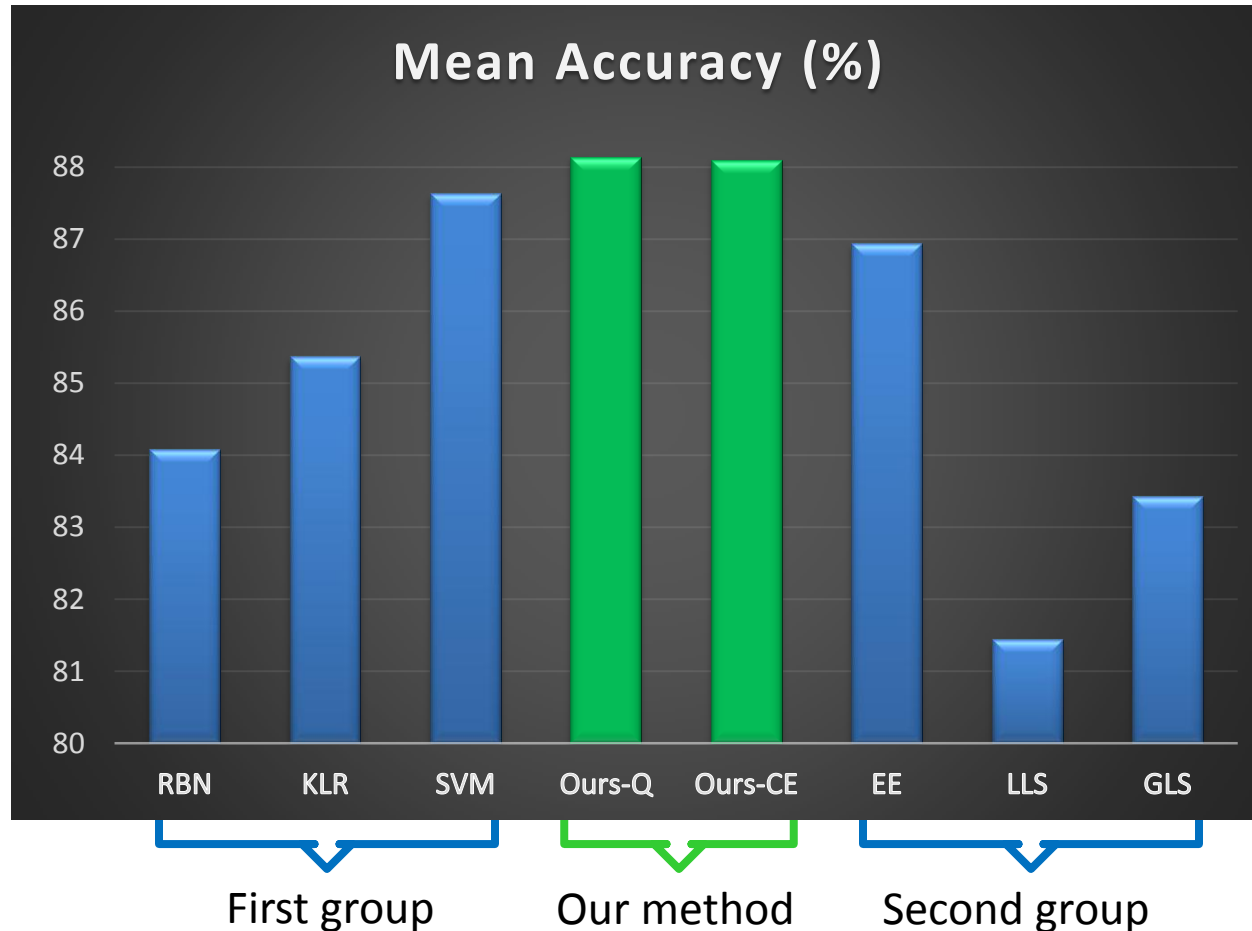
GLS vs. Ours-CE/Ours-Q

- same: RBF-based, volume based geometric regularizer
- diff geometry: on decision boundary vs on class probability

EE vs. Ours-Q

- same: RBF-based, quadratic loss
- diff geometric regularizer:
sophisticated on decision boundary vs on class probability

Results on UCI Datasets



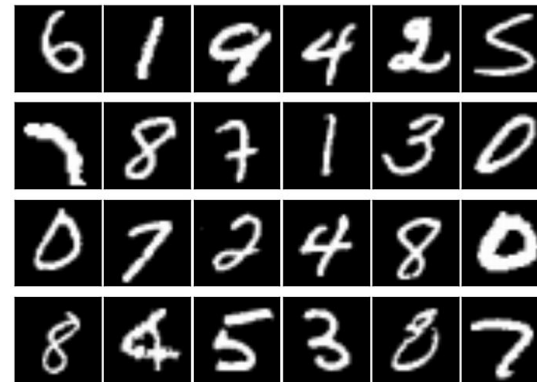
Experiments – RBF Representation

Real-world datasets – comparing with baseline

- Flickr Material Database (4096 dimensional feature)
- MNIST handwritten digits (60,000 samples)



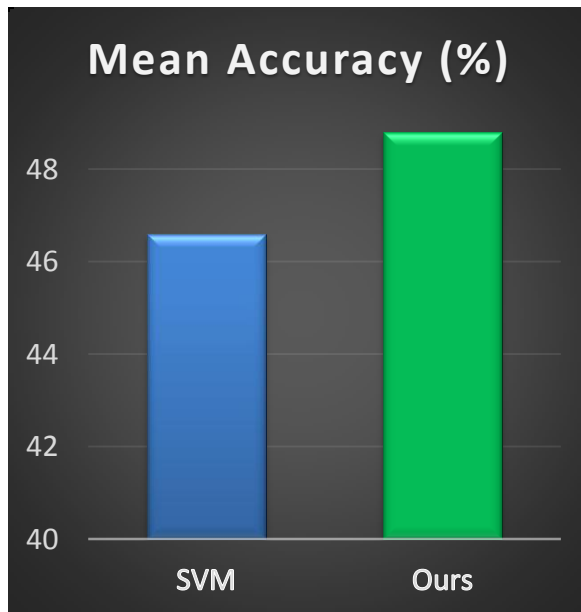
Flickr Material Database



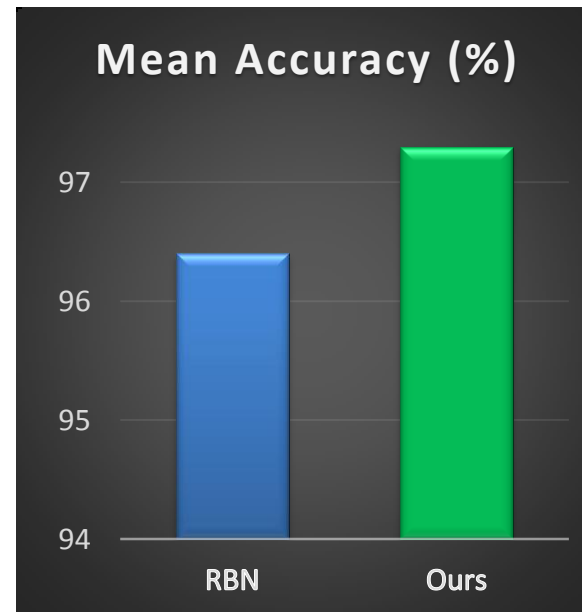
MNIST handwritten digits

Results on Real-world Datasets

Flickr Material Database



MNIST handwritten digits



Summary

- New geometric perspective on overfitting
- First regularization approach that exploits the geometry of a class probability estimator for classification
- Unified framework for both binary and multiclass cases
- Compares favorably to existing regularization methods

Collaborators



Steven Rosenberg
Mathematics, BU



Stan Sclaroff
Computer Science, BU



Zheng Wu
Mathworks Inc.