Machine Learning on Functional Data

Barnabás Póczos

June 9, 2015



MACHINE LEARNING DEPARTMENT





Motivation

- □ Function Classification/Regression
- Finding anomalous functions and distributions
- □ Function-to-Function Regression
- Fusso = Functional Shrinkage and Selection Operator (Functional Lasso)

ML on complex objects

Most ML algorithms operate on feature vectors...



Questions:

- □ How to learn?
- □ What can we learn?
- □ How fast can we learn?
- □ How many sample points do we need?
- □ Supervised/unsupervised/active learning ...

Functional Data is Everywhere















Motivating Examples

Manchester United 07/08







Owen Hargreaves Rio Ferdinand Cristiano Ronaldo

Shot Type

Goals
 Shots on Goal

Shots

www.juhokim.com/projects.php



Distribution Regression / Classification



ML on Distributions



SMAdardensadhineblaarning



What happens if we repeat the medical tests?

ML on Distributions



Dealing with complex objects

□ break into smaller parts, represent the input as a set of smaller parts

- □ treat the set elements as sample points from some **unknown distribution**
- □ do *ML on these unknown distributions* represented by sets

Distribution Classification

We have T sample sets, $(\mathbf{X}_1, \ldots, \mathbf{X}_T)$. [Training data] $\{X_{t,1}, \ldots, X_{t,m_t}\} = \mathbf{X}_t \sim p_t$. \mathbf{X}_t has class $Y_t \in \{-1, +1\}$.

What is the class label Y of $\mathbf{X} = \{X_1, \ldots, X_m\} \sim p$?

Solution: Use RKHS based SVM!

Calculate the Gram matrix $K_{ij} \doteq \langle \phi(p_i), \phi(p_j) \rangle_{\mathcal{K}} = K(p_i, p_j)$

Dual form of SVM: $\hat{\alpha} = \arg \max_{\substack{\alpha \in \mathbb{R}^T \\ T}} \sum_{i=1}^T \alpha_i - \frac{1}{2} \sum_{i,j}^T \alpha_i \alpha_j y_i y_j K_{ij}, \quad \text{subject to } \sum_i \alpha_i y_i = 0,$ $Y = \operatorname{sign}(\sum_{i=1}^T \hat{\alpha}_i y_i K(p_i, p)) \in \{-1, +1\}$ $0 \le \alpha_i \le C.$

Problems: We do not know p_i , p, $K(p_i, p_j)$, or $K(p_i, p)$...

Object Classification ETH-80 [Leibe and Schiele, 2003]



8 categories, 400 images, each image is represented by 576 18 dim points



Póczos, Xiong, Sutherland, & Schneider, CVPR 2012

Outdoor Scenes Classification [Oliva and Torralba, 2001]



Póczos, Xiong, Sutherland, & Schneider, CVPR 2012 11

Sport Events Classification [Li and Fei Fei, 2007]



badminton bocce croquet polo climbing rowing sailing snowboard 8 categories, 1040 images, each represented by 295 to 1542 57 dim points.



□ Best published: 86.7% (Zhang et al, CVPR 2011)

□ NPR: **87.1**%

Póczos, Xiong, Sutherland, & Schneider, CVPR 2012 12

Emerging images



Niloy J. Mitra, 2009

Emerging images



Image Representation with Distributions

Dealing with complex objects

- break into smaller parts,
- represent the object as a sample set of these parts



Image patches •Overlapping •Non-overlapping

Patch locations

- Grid pointsInteresting pointsRandom
- Patch sizes •Same •Different, •Hierarchy

d-dimensional sample set representation of the image

- □ Each image *patch* is represented by PCA compressed SIFT vectors.
 SIFT = *Scale-invariant feature transform*. *PCA: 128dim⇒ d dim* □ Each *image* is represented as a *set* of these *d* dim feature vectors.
- □ Each *set* is considered as a sample set from some *unknown distribution*.

Detecting Anomalous Images

B. Póczos, L. Xiong & J. Schneider, UAI, 2011.



2-dimensional sample set representation of images (128 dim SIFT \Rightarrow 2 dim) **Anomaly score:** divergences between the distributions of these sample sets

Detecting Anomalous Images

















51 52 53 54 55

Noisy USPS Dataset Classification with SDM



Original (noiseless) USPS dataset is easy ~97%
 Each instance (image) is a set of 500 2d points
 1000 training and 1000 test instances
 160



Results:

- SVM on raw images $82.1 \pm .5\%$ accuracy
- **SDM** on the 2D distributions, Rényi divergence: $96.0 \pm .3\%$ accuracy

Multidimensional Scaling of USPS Data

10 instances from figures 1,2,3,4.Calculate pairwise Euclidean distances.Nonlinear embedding with MDS into 2d.





Raw **images** using Euclidean distance



Estimated Euclidean distance between the **distributions**

Local Linear Embedding of Distributions

72 rotated COIL froggies





Edge detected COIL froggy







Finding Unusual Galaxy Clusters



What are the most anomalous galaxy clusters?

The most anomalous galaxy cluster contains mostly

- □ star forming blue galaxies
- □ irregular galaxies

B. Póczos, L. Xiong & J. Schneider, UAI, 2011. Credits: ESA, NASA

Understanding Turbulences







Turbulence Data Classification



Finding Vortices



Classification probabilities

Find Interesting Phenomena in Turbulence Data

Anomaly detection with 1-class SDM



Anomaly scores

Function-to-Function regression

Function to Function (F2F) Regression

Given an *input function p,* can we predict an *output function q*?

Co-occurring Prediction



Image Segmentation



Future Distribution Prediction

Data-Model

We observe functional observation *pairs* (P_i , Q_i), for i=1,...,N corresponding to input/output functions (p_i , q_i), where f(p_i) = q_i .

We predict output function \mathbf{q}_0 given unseen observation \mathbf{P}_0 .



$$P_{i} = \left\{ p_{i}(u_{ij}) + \epsilon_{ij} \right\}_{j=1}^{n_{i}}, Q_{i} = \left\{ q_{i}(v_{ij}) + \xi_{ij} \right\}_{j=1}^{m_{i}},$$

Kernel Methods Don't Scale

Previous Approaches:

- **RKHS** approach (Kadri et al. 2010)
- □ Kernel smoothing approach (Oliva et al. 2013)

$$(K(\tilde{p}_1, \tilde{p}_0), \dots, K(\tilde{p}_N, \tilde{p}_0))^T$$

N evaluations

$$\begin{pmatrix} K(\tilde{p}_{1}, \tilde{p}_{1}) & \cdots & K(\tilde{p}_{1}, \tilde{p}_{N}) \\ \vdots & \ddots & \vdots \\ K(\tilde{p}_{N}, \tilde{p}_{1}) & \cdots & K(\tilde{p}_{N}, \tilde{p}_{N}) \end{pmatrix}$$

Kernel Methods Don't Scale

Kernel smoothing approach

 $\widehat{f}(\widetilde{p}_0) = \sum_{i=1}^N W(\widetilde{p}_i, \widetilde{p}_0)\widetilde{q}_i$, where

$$W(\tilde{p}_i, \tilde{p}_0) = \begin{cases} \frac{K(D(\tilde{p}_i, \tilde{p}_0))}{\sum_{j=1}^N K(D(\tilde{p}_j, \tilde{p}_0))} & \text{if } \sum_{j=1}^N K(D(\tilde{p}_j, \tilde{p}_0)) > 0\\ 0 & \text{otherwise }. \end{cases}$$

D(p,q), is a distance or divergence between densities p and q.

Our Approach: Triple Basis Estimator (3BE)

High-level idea:

Embed input functions p with a *non-linear* map z(p)
 then use a *linear* map to estimate q.

$p \longrightarrow z(p) \longrightarrow q$

Orthonormal Basis

Let
$$\{\varphi_i\}_{i\in\mathbb{Z}}$$
 be an orthonormal basis for $L_2([0,1])$.
 $\{\varphi_\alpha\}_{\alpha\in\mathbb{Z}^d}$ where $\varphi_\alpha(x) = \prod_{i=1}^d \varphi_{\alpha_i}(x_i), x \in [0,1]^d$.

Let
$$h \in L_2([0,1]^d)$$
, then
 $h(x) = \sum_{\alpha \in \mathbb{Z}^d} a_\alpha(h) \varphi_\alpha(x)$
where $a_\alpha(h) = \langle \varphi_\alpha, h \rangle$
 $= \int_{[0,1]^d} \varphi_\alpha(z) h(z) dz \in \mathbb{R}.$



Orthogonal Projection Embedding

Calculating a(p) projection coefficient in the training set (Basis 1):



Function from Orthogonal Projections

Calculating a(q) projection coefficients in the training set (Basis 2):



$$\tilde{q}_j(x) = \sum_{k=1} \langle \tilde{q}_j, \varphi_{\beta_k} \rangle \varphi_{\beta_k}(x)$$

 $\mathbf{p} \longrightarrow \mathbf{a}(\mathbf{p}) \longrightarrow \mathbf{z}(\mathbf{a}(\mathbf{p})) \longrightarrow \mathbf{a}(\mathbf{q}) \longrightarrow \mathbf{q}_{34}$

Function-to-Function mapping

$$\tilde{p}_{j} \longrightarrow \tilde{q}_{j} \longrightarrow \tilde{q} \longrightarrow \tilde{q}_{j} \longrightarrow \tilde{q}_{j} \longrightarrow \tilde{q}_{j} \longrightarrow \tilde{q}_{j} \longrightarrow \tilde{q}_{j} \longrightarrow \tilde{$$

r regression problems (Multi task learning) $p \implies a(p) \implies z(a(p)) \implies a(q) \implies q_{a}$

Function-to-Function mapping

The kth regression problem:

Map \tilde{p}_j , (or its projection coefficients $\vec{a}(\tilde{p}_j)$) to the kth projection coefficent of \tilde{q}_j , $a_{\beta_k}(\tilde{q}_j)$ (for all j = 1, ..., N).

$$< \mathbf{M}, \mathbf{N} > = a_{\beta_k}(\tilde{q}_j) \approx \hat{f}_{\beta_k}(\tilde{p}_j) \qquad \vec{a}(\tilde{p}_j) = \begin{bmatrix} a_{\alpha_1}(\tilde{p}_j) \\ \vdots \\ a_{\alpha_s}(\tilde{p}_j) \end{bmatrix}$$

Let us use a kernel smoother for this regression:

$$\hat{f}_{\beta_k}(\tilde{p}_j) = \sum_{i=1}^N \theta_i K_{\sigma}(\|\tilde{p}_i - \tilde{p}_j\|_2) = \sum_{i=1}^N \theta_i K_{\sigma}(\|\vec{a}(\tilde{p}_i) - \vec{a}(\tilde{p}_j)\|_2) \qquad K_{\sigma}(\|\tilde{p}_i - \tilde{p}_j\|_2) = \exp\left(-\frac{\|\tilde{p}_i - \tilde{p}_j\|_2}{2\sigma^2}\right)$$

Random Kitchen Sinks

Theorem [Rahimi & Recht, NIPS 2007]

If one has a shift-invariant kernel K(in particular we consider the RBF kernel $K(x) = \exp(-x^2/2)$), then for fixed $\omega_i \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^{-2}I_d)$, $b_i \stackrel{iid}{\sim} U([0, 2\pi])$, we have that for each $x, y \in \mathbb{R}^d$:

$$K(\|x - y\|_2 / \sigma) \approx z(x)^T z(y),$$

where $z(x) \equiv \sqrt{\frac{2}{D}} \left[\cos(\omega_1^T x + b_1) \cdots \cos(\omega_D^T x + b_D) \right]^T,$

and D is the number of random basis functions (see Rahimi & Recht (NIPS 2007) for approximation quality).

Fast Projection Estimation (Basis 3)

$$< M > = a_{\beta_k}(\tilde{q}_j) \approx \hat{f}_{\beta_k}(\tilde{p}_j)$$

$$\hat{f}_{\beta_k}(\tilde{p}_0) = \sum_{i=1}^N \theta_i K_\sigma(\|\tilde{p}_i - \tilde{p}_0\|_2)$$

$$= \sum_{i=1}^N \theta_i K_\sigma(\|\vec{a}(\tilde{p}_i) - \vec{a}(\tilde{p}_0)\|_2)$$

$$\approx \sum_{i=1}^N \theta_i z(\vec{a}(\tilde{p}_i))^T z(\vec{a}(\tilde{p}_0))$$

$$= \left(\sum_{i=1}^N \theta_i z(\vec{a}(\tilde{p}_i))\right)^T z(\vec{a}(\tilde{p}_0))$$

$$K_{\sigma}(\|\tilde{p}_i - \tilde{p}_j\|_2) = \exp\left(-\frac{\|\tilde{p}_i - \tilde{p}_j\|_2^2}{2\sigma^2}\right)$$

RKS (Rahimi & Recht 2007):

$$K_{\sigma}(\|\vec{a}(\tilde{p}_{i}) - \vec{a}(\tilde{p}_{j})\|_{2}) \approx z(\vec{a}(\tilde{p}_{i}))^{T} z(\vec{a}(\tilde{p}_{j}))$$
$$z(\vec{a}(\tilde{p}_{j})) = \begin{bmatrix} z_{1}(\vec{a}(\tilde{p}_{j})) \\ \vdots \\ z_{D}(\vec{a}(\tilde{p}_{j})) \end{bmatrix}$$

 $= \psi_{\beta_k}^T z(\vec{a}(\tilde{p}_0)) \qquad \text{D dim linear regression!}$ We have r of these. This is the kth.

where $\psi_{\beta_k} = \sum_{i=1}^N \theta_i z(\vec{a}(\tilde{p}_i)) \in \mathbb{R}^D$.

$p \Longrightarrow a(p) \longrightarrow z(a(p)) \longrightarrow a(q) \longrightarrow q_{38}$

Sample Projection Function Estimate

Let us write the r D-dim linear regression in a matrix form on the N training data

 $\vec{a}(\tilde{q}_0) \approx \hat{\Psi}^T z(\vec{a}(\tilde{p}_0))$ where



$$\hat{\Psi} \equiv \underset{\Psi \in \mathbb{R}^{D \times r}}{\arg\min} \|A - \mathbf{Z}\Psi\|_F^2 = (\mathbf{Z}^T \mathbf{Z})^{-1} \mathbf{Z}^T A$$

 $p \longrightarrow a(p) \longrightarrow z(a(p)) \longrightarrow a(q) \longrightarrow q_{39}$

Computational Complexity for Prediction

- N : # of instances (input/output pairs)
- n : # of sample points per function observation
- **S** : # of projection coefficients for p
- *r* : # of projection coefficients for q
- D: # of random basis functions

The total computation time for prediction is O(rD + Ds + sn).

If $D = O(n \log(n))$, s = O(n), and r = O(n), then, the total runtime for evaluating $\hat{f}(\tilde{p}_0)$ is $O(n^2 \log(n))$.

Since we are considering data-sets where the number of instances N far outnumbers the number of points per sample set $n, O(n^2 \log(n))$ is a substantial improvement over $\Omega(Nn)$ for the LSE.

Risk Upper Bound

$$\mathbb{E}\left[\|q_0 - \hat{q}_0\|_2^2\right] \le O\left(\left(n^{-1/(2+\gamma_{\mathcal{I}}^{-1})} + \frac{n\log(n)\log(N)}{N}\right)^{2/(2+\gamma_{\mathcal{O}}^{-1})}\right)$$

- N : # of instances (input/output pairs)
- n : # of sample points per function observation
- $\gamma_{\mathcal{I}}^{-1}$: input function smoothness parameter
- $\gamma_{\mathcal{O}}^{-1}$: output function smoothness parameter

Rectifying 2LPT Simulations (2nd order Lagrangian Perturbation Theory) 2LPT z=0.125 z=0.167 z=0.208



ひ N-body



Method	MSE	MPT
3BE	4.958	0.009
LSE	6.816	4.977
2LPT	6.424	NA



Co-occurring Prediction with MoCap



Forward Prediction with Music Data



F2F Learning - Summary

Scalable Function-to-Function regression with Triple Basis Estimator (3BE)

- Orthonormal basis to represent input functions
- Orthonormal basis to represent output functions
- Random basis mapping input functions to output functions

Fusso = Functional Shrinkage and Selection Operator (Functional Lasso)

FuSSO

We present the FuSSO, a functional analogue to the LASSO:

- Finds a sparse set of functional input covariates to regress a real-valued response
- The FuSSO is semi-parametric: No parametric assumptions on input functional covariates
- □ Assumes linear form to response from functional covariates
- We provide a statistical backing for use of the FuSSO via proof of asymptotic sparsistency under various conditions
- Furthermore, we observe good results on both synthetic and real-world data

Real Valued Regression with Multiple Functional Covariates

Function to real regression is a previously studied problem (e.g. in functional analysis); here the mapping takes in one function:



Similarly, one may consider a mapping that takes in multiple functions:



The number of functional input covariates may be very large; thus, a sparse model that depends only on a few of the functional covariates may be preferred:



Example Applications

- Finance:
 - *Inputs:* Time-series of several commodities prices in the past
 - Output: Price of a particular commodity in the nearby future
- Neuroscience:
 - *Inputs:* Functions at each voxel (e.g. orientation distribution functions)
 - Output: The age of the subject





Linear Functional Regression

There are functional analogues to familiar real-vector linear regression models; for $Y_i \in \mathbb{R}$, $\epsilon_i \sim \mathcal{N}(0, \sigma)$, and $\Psi \subseteq \mathbb{R}^k$, a compact set:

One Real Vector vs. Functional Covariate: **Real Vector Covariate** Functional Covariate $Y_i = \langle X_i, w \rangle + \epsilon_i \mid Y_i = \langle f^{(i)}, g \rangle + \epsilon_i$ where $X_i, w \in \mathbb{R}^d$ and $\langle X_i, w \rangle = \sum_{j=1}^d X_{ij} w_j \mid f^{(i)}, g \in L_2(\Psi)$ and $\langle f^{(i)}, g \rangle = \int_{\Psi} f^{(i)}(t)g(t)dt$ Similarly, instead of having one input covariate per data instance, one may have p feature vectors, $\{X_{i1}, \ldots, X_{ip}|, X_{ij} \in \mathbb{R}^d\}$, or functions, $\{f_1^{(i)}, \ldots, f_p^{(i)} \mid f_j^{(i)} \in L_2(\Psi)\}$, associated to each data instance i:

Multiple Real Vectors vs. Functional Covariates:

 $\begin{array}{l|l} \mbox{Real Vector} & \mbox{Functional Covariates} \\ Y_i = \sum_{j=1}^p \langle X_{ij}, w_j \rangle & Y_i = \sum_{j=1}^p \langle f_j^{(i)}, g_j \rangle \\ & y_i = \sum_{j=1}^p \langle f_j^{(i)}, g_j \rangle \\ & where & + \epsilon_i \\ & + \epsilon_i \\ & w_1, \dots, w_p \in \mathbb{R}^d & g_1, \dots, g_p \in L_2(\Psi) \end{array}$

FuSSO: Functional Shrinkage and Selection Operator

Regression with Basis Functions:

Let $\Psi = [0, 1], \Phi = \{\varphi_m\}_{m=1}^{\infty}$ be an orthonormal basis for $L_2(\Psi)$:

 $f_j^{(i)}(x) = \sum_{m=1}^{\infty} \alpha_{jm}^{(i)} \varphi_m(x), \text{ where } \alpha_{jm}^{(i)} = \int_0^1 f_j^{(i)}(t) \varphi_m(t) dt$

Similarly $g_j(x) = \sum_{m=1}^{\infty} \beta_{jm}^* \varphi_m(x)$, then by orthonormality:

$$\langle f_j^{(i)}, g_k \rangle = \sum_{m=1}^{\infty} \alpha_{jm}^{(i)} \beta_{km}^*$$

Let $\{\tilde{\alpha}_{jm}\}_{m=1}^{M_n}$ approximate $\{\alpha_{jm}\}_{m=1}^{M_n}$, and
 $f_j^{(i)}(x) \approx \tilde{f}_j^{(i)}(x) = \sum_{m=1}^{M_n} \tilde{\alpha}_{jm}^{(i)} \varphi_m(x)$

Optimization Problem:

We may then estimate $g_j = \sum_{m=1}^{\infty} \beta_{jm}^* \varphi_m$ as $\hat{g}_j = \sum_{m=1}^{M_n} \hat{\beta}_{jm} \varphi_m$ where $\hat{\beta} \equiv \operatorname{argmin}_{\beta} \frac{1}{2N} \sum_{i=1}^{N} \left(Y_i - \sum_{j=1}^{p} \sum_{m=1}^{M_n} \tilde{\alpha}_{jm}^{(i)} \beta_{jm} \right)^2$ $+ \lambda_N \sum_{j=1}^{p} \sqrt{\sum_{m=1}^{M_n} \beta_{jm}^2}$

Theory: Sparsistency

 $\mathbb{P}\left(\operatorname{supp}(\hat{\beta}) = \operatorname{supp}(\beta^*)\right) \to 1 \text{ when given a}$ dataset $\mathcal{D} = \{(\{\vec{y}_j^{(i)}\}_{j=1}^p, Y_i)\}_{i=1}^N, \text{ where}$

$$\begin{split} \vec{y}_{j}^{(i)} &= \vec{f}_{j}^{(i)} + \xi_{j}^{(i)}, \\ \vec{f}_{j}^{(i)} &= \left(f_{j}^{(i)}(1/n), \ f_{j}^{(i)}(2/n), \ \dots, \ f_{j}^{(i)}(1) \right)^{T}, \\ \xi_{j}^{(i)} &\stackrel{iid}{\sim} \mathcal{N}(0, \sigma_{\xi}^{2} I_{n}), \end{split}$$

and $f_j^{(i)}, g_j : [0, 1] \mapsto \mathbb{R}$. The proof is generalizable to higher dimensions, and when $f_j^{(i)}$ are pdfs and one observes samples.

Results Synthetic Dataset:

• Tested on datasets of $\mathcal{D} = \{(\{\vec{y}_j^{(i)}\}_{j=1}^p, Y_i)\}_{i=1}^N$

 $-\vec{y}_{j}^{(i)}$ grid of *n* noisy grid function evaluations

- The number of true functional covariates in the support, s, was fixed to be 5
- Functions $f_j^{(i)}$ and g_j were generated randomly

For each experiment we ran 100 trials and recorded the following:

Experiments' Variables:

Name	Definition		
p	# of input covariates		
N	# of data instances		
n	# of grid evaluations		
r	% of trails w/ support found		
Δ	avg. range of λ w/ support		

Experiments' Results:(p, N, n)r Δ (100, 50, 5).68.2125(1000, 500, 25)1.4771(20000, 500, 25)1.4729



Neurological Dataset:

- Dataset with over 25K functions per subject for 89 total subjects
 - Orientation distribution functions (ODF) at white matter voxels
- We regress a subject's age, given ODFs
- We compared to using the LASSO with peak ODF (quantitative anisotropy) values

Results					
Method:	FuSSO (ODFs)	LASSO (QAs)	Mean Predict		
MSE:	70.85	77.13	156.43		





mage Sources: http://www.aging2.com/wp-content/uploads/2013/05/Screen-Shot-2013-05-28-at-9.48.49-PM.png;

http://media.salon.com/2013/02/money1.jpg; http://3278as3udzze1hdk0f2th5nf18c1.wpengine.netdna-cdn.com/wp-content/uploads/2010/10/connectome-brain-diffusion-spectrum-imaging.j62

If you are interested, contact me! © **bapoczos@cs.cmu.edu, GHC-8231**

Functional data has so many applications!

Some results on regression/classification/anomaly detection/ Lasso

Lots of missing theoretical results: Lower bounds, active learning



Thanks for your attention! ③

