Machine Learning and Pattern Recognition

Sargur N. Srihari University at Buffalo, The State University of New York USA

Adobe Distinguished Lecture, Bangalore, India December 11, 2009

National Academy of Sciences: Keck Center



Regression Problem: Carbon Dioxide in Atmosphere



Classification Problem: Handwritten Digit Recognition



Wide variability of same numeral

- Handcrafted rules will result in large no of rules and exceptions
- Better to have a machine that learns from a large training set

Types of Problems ML is used for

Supervised Learning Classification 0123456789

Regression



Unsupervised Learning

Clustering



$$p(\mathbf{x}) = \sum_{k=1}^{K} \pi_k N(\mathbf{x} \mid \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

5

Regression Problem: Robot Kinematics



Input: Link Angles

Output: Effector Position

Regression Problem: Learning to Rank



	(Geogle Beanth) (Tri Feeling Lucky)	
	1998 Images News Miles	
bing	1	
<u> </u>	E three al. Conty from Makyala	

Input: (Query-URL) pair

Output: Relevance Value









Regression Problem: Biometrics/Forensics



Input: Minutiae

Output: Core Point

The ML Approach



(Inference

OR

Testing)

Search

Find optimal solution to problem

ML History

- ML has origins in Computer Science
- PR has origins in Engineering
- They are different facets of the same field



- Recent Advances
 - Fully Bayesian methods
 - availability of computational methods
 - New Models
 - Gaussian Processes, CRF



Issues in Regression Problems

- Robot Kinematics
 - Forward/inverse problems
- Ranking Web Pages
 - High-dimensional input vectors
- Carbon dioxide Emission
 - Periodic data
- Fingerprint Analysis
 - Multi-dimensional output

Kinematics of a Robot Arm

Forward problem: Find end effector position given joint angles Has a unique solution

Forward problem:Inverse kinematics: two solutions:Find end effector positionElbow-up and elbow-down





 Forward problems correspond to causality in a physical system have a unique solution

e.g., symptoms caused by disease

• If forward problem is a many-to-one mapping, inverse has multiple ¹²/_{solutions}

Robot Kinematics: Forward/Inverse Problems

Forward data set



Red curve obtained by minimizing sum-of-squared error Inverse data set: Reversing *x* and *t*



Very poor fit to data: Mixture Density Network



Search Engine Ranking: Features



LEarning TO Rank In LETOR 4.0 dataset 46 query-document features Maximum of 124 URLs (For one query)

Features (Google uses over 200)

- Log frequency of query word in anchor text
- Query word in color on pags
- # of images on page
- # of (out) links on page
- PageRank of page
- URL length
- URL contains "~"
- Page length

Query used in some features

Traditional IR used few document features

14

- Term Frequency
- Inverse Document Frequency

Target Variable

- Point-wise (0,1,2,3)

Statistical Models

Generative

- Naïve Bayes
- Mixtures of multinomials
- Mixtures of Gaussians
- Hidden Markov Models
 (HMM)
- Bayesian networks
- Markov random fields

- Discriminative
 - Logistic regression
 - SVMs
 - <u>Traditional neural</u>
 <u>networks</u>
 - Nearest neighbor
 - Conditional Random Fields (CRF)
 - Gaussian Processes

Linear Models

- Regression
 - Single input variable
 - Polynomial $y(x,w) = w_0 + w_1 x + w_2 x^2 + ... = \sum w_i x^i$
 - Several variables
 - linear combination of non-linear (basis) functions

$$y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{j=1}^{M-1} w_j \phi_j(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x})$$

– Bayesian Linear Regression



Classification

- Logistic Regression (with sigmoid or soft-max)

 $y(\mathbf{x},\mathbf{w}) = \sigma[\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x})]$

Neural Network Function

Overall function

$$y_{k}(\mathbf{x},\mathbf{w}) = \sigma \left(\sum_{j=1}^{M} w_{kj}^{(2)} h \left(\sum_{i=1}^{D} w_{ji}^{(1)} x_{i} + w_{j0}^{(1)} \right) + w_{k0}^{(2)} \right)$$

- Where w is the set of all weights and bias parameters
 - nonlinear functions from inputs $\{x_i\}$ to outputs $\{y_k\}$
 - controlled by vector w of adjustable parameters
 - Note presence of both σ and h functions
- If σ is identity for regression
- If σ is sigmoid for two-class classification
- If σ is softmax for multi-classification



Gaussian Processes for Regression

- Radically different viewpoint
 - not involving weight parameters
- Functions are drawn from a Gaussian
 - where each data point is a function



Gaussian Process Regression

- Generalize multivariate Gaussian to infinite variables
- A Gaussian *distribution* is fully specified by a mean vector μ and covariance matrix Σ $f = (f_1, ..., f_n)^T \sim \mathcal{N}(\mu, \Sigma)$ indexes i = 1, ..., n
- A Gaussian process is fully specified by a mean function m(x) and covariance function k(x,x')
 - $f(x) \sim GP(m(x), k(x, x'))$ indexes x
 - Kernel function *k* appears in place of covariance matrix
 - Both express similarity of two in multidimensional space

Gaussian Process with Two Samples

- Let y be a function (curve)
 of a one-dimensional variable x
- We take two samples y_1 and y_2 corresponding to x_1 and x_2
- Assume they have a bivariate Gaussian distribution
- Each point from this distribution
 - has an associated probability
 - It also defines a function y(x)
 - Assuming that two points are enough to define a curve
- More than two points will be needed to define a curve
 - Which leads to a higher dimensional probability distribution



The marginalization property

- GP as a Gaussian with an infinitely long mean vector and an infinite by infinite covariance matrix may seem impractical
- Luckily we have the marginalization property $p(x) = \int p(x,y) dy$
 - For Gaussians

$$p(\mathbf{x},\mathbf{y}) = N \begin{pmatrix} \begin{bmatrix} a \\ b \end{bmatrix}, \begin{bmatrix} A & B \\ B^T & C \end{bmatrix} \implies p(\mathbf{x}) = N(\mathbf{a},A)$$

- Marginalize over unknown function values 1.
- To get distribution over known values x





Predictive distribution

- Mean $m(\mathbf{x}_{N+1}) = \mathbf{k}^{\mathrm{T}} \mathbf{C}_{N}^{-1} \mathbf{t}$
- Variance σ²(x_{N+1}) = c k^TC_N⁻¹k where C_N is the N x N covariance matrix with elements k(x_n,x_m) + β⁻¹δ_{nm} vector k has elements k(x_n,x_{N+1}) for n=1,.., N and scalar c = k(x_{N+1},x_{N+1})+β⁻¹
- Restrictions on kernel function: positive definite



Periodic Smooth Functions

To create distribution over periodic functions of *x*:

- First map inputs to $u = (\sin x, \cos x)^T$
- Then measure distances in *u*-space
- Combined with squared exponential (SE) covariance function with characteristic length ale **t** we get



²⁴ Three functions drawn at random: left l > 1, right l < 1

Covariance Function for Periodic Data

The covariance function consists of several terms, parameterized by 11 hyperparameters

$$k(x,x') = k_1(x,x') + k_2(x,x') + k_3(x,x') + k_4(x,x')$$

Where

Long-term smooth trend (squared exponential) $k_1(x,x')=q_1^2\exp(-(x-x')^2/q_2^2)$ Seasonal trend (quasi-periodic smooth) $k_2(x,x')=q_3^2\exp(-2sin^2[p(x-x')]/q_s^2 \times exp(-(1/2)(x-x')^2/q_4^2)),$ Short-term and medium-term anomaly (rational quadratic) $k_2(x,x')=q_2^2(1+\frac{(x-x')^2}{2})=0$

$$k_{3}(x,x') = \theta_{6}^{2} \left(1 + \frac{(x-x')^{2}}{2\theta_{8}\theta_{7}} \right) - \theta_{8}$$

Noise (independent Gaussian and dependent)

$$k_4(x,x') = \theta_9^2 \exp\left(-\frac{(x-x')2}{2\theta_{10}}\right) + \theta_{11}^2 \delta_{xx'}$$

25

Computation of Fingerprint Probabilities







Madrid Latent

Brandon Mayfield

Daoud

Core Point Identification for Latent Fingerprints

- Goal:
 - Given a latent (partial) fingerprint, predict its core point
 - Core point is the center of the north most loop type singularity
 - Latent fingerprints are represented by orientation maps.



Latent fingerprint and its matching 10-print mate.

Orientation Map

Cross: core point predicted in latent print

Circle: True core point in corresponding 10-print

Regression by Gaussian Processes

- Training set of *N* fingerprints: $D = \{(g_i, y_i) | i = 1, ..., N\}$

 - -g = orientation map and y = core point
- Regression model with Gaussian noise •

$$y = f(g) + \varepsilon$$
 $\varepsilon \sim N(\varepsilon \mid 0, \delta^2)$

- $-\epsilon$ is a random noise variable:
- Likelihood for core point: •

$$p(y | f(g)) = N(f(g), \delta^2 I)$$

Gaussian process prior:

 $f(g) \sim GP(0, k(g, g'))$

- where k(g, g') is the kernel function used to define the covariance, e.g.

$$k(g,g') = \theta_1 \exp(-\theta_2 |g-g'|^2)$$

Predictive Distribution

- Let $g = (g_1, ..., g_N)$ and $y = (y_1, ..., y_N)$
 - Orientation map of latent print is g^*
 - Predictive distribution is given by $p(y^* | g^*, g, y) = N(m(y^*), \operatorname{cov}(y^*))$ $m(y^*) = k(g^*, g)[K + \delta^2 I]^{-1}y$

- where

$$\operatorname{cov}(y^*) = k(g^*, g^*) + \delta^2 - k(g^*, g)^T [K + \delta^2 I]^{-1} k(g, g^*)$$

• *K* is the Gram matrix whose elements are given by $k(g_i, g_j)$

Core Point Learning Results Latent cases from NIST data set Case "b115" Case "g73"



Table 1. Comparison of prediction precisions of Poincare index and Gaussian processes based approaches.

	Poincare Index	Gaussian Processes
Good	90.6%	93.1%
Bad	68.2%	87.1%
Ugly	46.6%	72.7%
Overall	68.6%	84.5%

Probability of Random Correspondence

After core points are aligned, feature vectors are computed to determine probability



Summary of ML

- Principled way of building high performance information processing systems
 - -Systematic approach instead of ad-hockery
- Useful for solving problems with inherent uncertainty
 - Supervised Learning
 - Classification, Regression
 - Unsupervised Learning
 - Clustering
- Wide range of applications

Summary of ML Applications

- Artificial Intelligence

- Perform tasks that humans perform well but difficult to specify algorithmically
- Information Retrieval
 - Search engines
 - Web personalization
- Scientific applications
 - Computational science
- Engineering

0123456789









NAS Committee Report





Released March 2009 National Academies Press

Committee at NAS, Woods Hole, MA