Machine Learning: New Paradigms and Perspectives

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

Outline

- 1. What is Machine Learning (ML)?
- 2. The Fully Bayesian Approach
 - Biometric/Forensic Application
- 3. Graphical Models
- 4. Discriminative and Generative Models
 - Conditional Random Fields
 - MINDS Applications
- 5. Concluding Remarks

Example 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

What is Machine Learning?

- Programming computers to:
 - Perform tasks that humans perform well but difficult to specify algorithmically (PR)
 - Automatically improve with experience (AI)
- Principled way of building high performance information processing systems
 - search engines (IR)
 - adaptive user interfaces, personalized assistants (information systems)
 - scientific application (computational science)

ML History

- ML has origins in Computer Science
- PR has origins in Engineering
- They are different facets of the same field
- Methods around for over 50 years
- Revival of Bayesian methods

 due to availability of computational methods

The ML Approach

Data Collection Samples of human performance Model Selection Probability distribution to model process Parameter Estimation Values/distributions





Search

Find optimal solution to problem

Decision (Inference OR Testing)

The Fully Bayesian Approach

- Bayes Rule
- Bayesian Probabilities
- Concept of Conjugacy
- Monte Carlo Sampling

Rules of Probability

• Given random variables X and Y

• Sum Rule gives Marginal Probability

$$p(X = x_i) = \sum_{j=1}^{L} p(X = x_i, Y = y_j) = \frac{c_i}{N}$$



• Product Rule: joint probability in terms of conditional and marginal n_{ii} n_{ii} n_{ii} c_i

$$p(X,Y) = \frac{n_{ij}}{N} = p(Y | X) p(X) = \frac{n_{ij}}{c_i} \times \frac{c_i}{N}$$

• Combining we get Bayes Rule

$$p(Y | X) = \frac{p(X | Y)p(Y)}{p(X)}$$
 where $p(X) = \sum_{Y} p(X | Y)p(Y)$

Viewed as

Posterior α likelihood x prior

Bayesian Representation of Uncertainty

- Not just frequency of repeatable event
- A quantification of uncertainty
- Example: Whether Arctic ice cap will disappear by 2100
 - We have an idea of how quickly polar ice is melting
 - Revise it on basis of fresh evidence (satellite observations)
 - Assessment affects actions we take (to reduce greenhouse gases)
- Handled by general Bayesian
 interpretation

- Use of probability to represent uncertainty is not an ad-hoc choice
- If numerical values represent degrees of belief,
 - then simple axioms for manipulating degrees of belief leads to sum and product rules of probability

Probability Distributions: Relationships



Fully Bayesian Approach

- Feasible with increased computational power
- Intractable posterior distribution handled using either
 - variational Bayes or
 - stochastic sampling
 - e.g., Markov Chain Monte Carlo, Gibbs

An application in biometrics/forensics

- Task involves significant uncertainty in some modalities
 - E.g., Face, handwriting

Verification Task Passport photo





Same Person?



Vivek 13

Verification Task





Same Person?

Harish

Person Independent Training

Same Person



Chen





Vivek









Harish





Vivek Harish

Different Person





Chen Pavitra





Harish Pavitra

Vivek

Chen





Person Dependent: Several Samples Available for Given Person

Training On One Person









Verification of Questioned



Questioned 1



Questioned 2

Genuine and Questioned Distributions



C





Genuine Distribution .4,.1,.2,.3,.2,.1 **Distribution** Comparison 0,.4,.1,.2 **Ouestioned** Distribution

Bayesian Comparison of Distributions

- S_1 Samples from one distribution
- S_2 Samples from another distribution
- Task: To determine statistical similarity between S_1 and S_2
- D_1, D_2 Theoretical distributions from which S_1 and S_2 came from
 - Determine Prob(D₁ = D₂) given that we only observe samples from them
 - Method:
 - Determine $Prob(D_1 = D_2 = a \text{ particular distribution } f)$
 - Marginalize over all f (different means and variances)
 - If *f* is Gaussian we obtain a closed-form solution for marginalization over the family of Gaussians F.

Bayesian Comparison Formulation

$$P(D_{1} = D_{2} = f \mid S_{1}, S_{2}) = \frac{p(S_{1}, S_{2} \mid f) p(f)}{\iint_{F} p(S_{1}, S_{2} \mid D_{1} = f_{1}, D_{2} = f_{2}) df_{1} df_{2}}$$

$$p(D_{1} = D_{2} \mid S_{1}, S_{2}) = \int_{F} \frac{p(S_{1}, S_{2} \mid f) p(f)}{\int_{F} p(S_{1} \mid F_{2} \mid f) p(f)}$$

Bayes rule

Marginalizing over all f $\int_{F} P(S_1 | D_1 = f_1) df_1 \times \int_{F} P(S_2 | D_2 = f_2) df_2 \quad F - Family of Gaussians$ variables of integration. If f is drawn **uniformly** from the family F, then p(f) = p(f1) = p(f2) = 1

$$p(D_1 = D_2 \mid S_1, S_2) = \frac{Q(S_1 \cup S_2)}{Q(S_1) \times Q(S_2)}, Q(S) = \int_F p(S \mid f) df$$

Evaluation of Q for Gaussian family

$$Q(S) = \int_{-\infty}^{\infty} \int_{0}^{\infty} P(S \mid \mu, \sigma) d\mu d\sigma$$
$$Q(S) = 2^{-3/2} \times \pi^{(1-n)/2} \times n^{(-1/2)} \times C^{(1-n/2)} \times \Gamma(n/2-1)$$
$$C = \sum_{x \in S} x^{2} - \frac{\left(\sum_{x \in S} x\right)^{2}}{n}$$

For Gaussians closed form for Q is obtainable

where *n* is the cardinality of S

Results – Error Rates



Graphical Models

What are they?

Marriage between probability theory and graph theory Advantage

Allows viewing many probabilistic systems as instances of a common underlying formalism

Directed/Undirected Graphical Models

- Directed
- Bayesian networks
 - Express causal relationships
- Al and statistics



- Undirected
- Markov random fields
 - express soft constraints
- Vision and physics



Two Models Described by Same Graph



- 1. Hidden Markov Model: If latent variables are discrete: Observed variables in a HMM may be discrete or continuous
- 2. Linear Dynamical Systems: If both latent and observed variables are Gaussian

Generative and Discriminative Methods

Generative Methods

- Sampling can generate synthetic data
- Discriminative Methods
 - Directly estimate posterior distributions
 - No attempt to model underlying probability distributions
 - Focus computational resources on given task- better performance

Popular ML Models

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

- Discriminative
 - Logistic regression
 - SVMs
 - Traditional neural networks
 - Nearest neighbor
 - Conditional Random Fields (CRF)

Relationships between Models Single decision vs Sequential



From HMM to CRF

• In an HMM

$$p(\mathbf{Y},\mathbf{X}) = \prod_{n=1}^{N} p(y_n | y_{n-1}) p(\mathbf{x}_n | y_n)$$

Can be rewritten as

$$p(\mathbf{Y}, \mathbf{X}) = \frac{1}{Z} \exp\left\{\sum_{m=1}^{M} \lambda_m f_m(y_n, y_{n-1}, \mathbf{x}_n)\right\}$$

• Which gives us

$$p(\mathbf{Y} \mid \mathbf{X}) = \frac{p(y, x)}{\sum_{y'} p(y', x)} = \frac{\exp\left\{\sum_{m=1}^{M} \lambda_m f_m(y_n, y_{n-1}, \mathbf{X}_n)\right\}}{\sum_{y'} \exp\left\{\sum_{m=1}^{M} \lambda_m f_m(y_n, y_{n-1}, \mathbf{X}_n)\right\}}$$

Functional Models



Application of ML to Language Technologies

- Information Retrieval (IR)
- Natural Language Processing (NLP)
- Document Analysis and Recognition (DAR)
- Automatic Speech Recognition (SR or ASR)

ML in IR

- IR is historically based on empirical considerations
 - Not concerned with whether based on theoretically sound principles
- Some IR Tasks where ML is used
 - Relevance Feedback
 - Use patterns of documents accessed in the past
 - Document Ranking (Separating wheat from chaff)
 - Using Server Logs
 - Document Gisting and Query Relevant Summarization
 - Using FAQ lists
 - Regularities in very large databases (Data Mining)

ML in NLP

- Part Of Speech tagging
- Table Extraction
- Shallow Parsing
- Named Entity Tagging
- Text Classification

NLP: Part Of Speech Tagging

For a sequence of words $w = \{w_1, w_2, ..., w_n\}$ find syntactic labels s for each word:

w = The quick brown fox jumped over the lazy dog s = DET VERB ADJ NOUN-S VERB-P PREP DET ADJ NOUN-S

Model	Error
HMM	5.69%
CRF	5.55%

Baseline is already 90%

✓ Tag every word with its most frequent tag

✓Tag unknown words as nouns

Per-word error rates for POS tagging on the Penn treebank

Table Extraction

- To label lines of text document:
 Whether part of table and its role in table
- A necessary component of data mining, question-answering and IR tasks.

HMM	CRF
89.7%	99.9%

Shallow Parsing

- Precursor to full parsing or information extraction
 - Identifies non-recursive cores of various phrase types in text
- Input: words in a sentence annotated automatically with POS tags
- Task: label each word with a label indicating

- word is outside a chunk (O), starts a chunk (B), continues a $NP^{Chunks}(I)$

Rockwell In	iterna	tional Corp.	's Tulsa un	it said	it	si	gned	a tentative agreem	ent	extending
its contract	with	Boeing Co.	to provide	structural	parts	for	Boeing	's 747 jetliners .		

CRFs beat all reported single-model NP chunking results on standard evaluation dataset

Model	F score
CRF	94.38%
Generalized winnow	93.89%
Voted perceptron	94.09%
MEMM	93.70%

CRF Model for Word Recognition

Probability of image \mathbf{x} is word '*the*' is:

$$P(\mathbf{y} = 'the' | \mathbf{x}, \theta) = \frac{e^{\psi(y = 't', x; \theta)}}{\sum_{y' = 'a'}^{'z'} e^{\psi(y', x; \theta)}} \frac{e^{\psi(y = 'h', y_n = 't', x; \theta)}}{\sum_{y' = 'a'}^{'z'} e^{\psi(y', y_n = 't', x; \theta)}} \frac{e^{\psi(y = 'e', y_n = 'h', x; \theta)}}{\sum_{y' = 'a'}^{'z'} e^{\psi(y', y_n = 't', x; \theta)}}$$

$$\psi(y, y_n \mathbf{x}; \theta) = A(j, y, \mathbf{x}; \theta^s) + I(j, k, y, y_n, \mathbf{x}; \theta^t)$$

Captures state features for a character

Height, width, aspect ratio, position in text, etc

Captures transition features between a character and its preceding character in the word

Vertical overlap Total width of the bigram Difference in the height, width, aspect ratio

CRF for Document Segmentation

$$P(\mathbf{y}|\mathbf{x},\theta) = \frac{e^{\psi(y,\mathbf{x};\theta)}}{\sum_{y'} e^{\psi(y',\mathbf{x};\theta)}}$$

where $\mathbf{y} \in \{\text{Machine-print}, \text{Handwriting}, \text{Noise}\}\$ $\theta : \text{CRF model parameters}\$ $\mathbf{x} : \text{Observed document}$

$$\psi(y, x; \theta) = \sum_{j=1}^{m} \left(A(j, y_j, \mathbf{x}; \theta^s) + \sum_{(j,k) \in E} I(j, k, y_j, y_k, \mathbf{x}; \theta^t) \right)$$

where $A(j, y_j, \mathbf{x}; \theta^s)$: association potential $I(j, k, y_j, y_k, \mathbf{x}; \theta^t)$: interaction potential m: non-overlapping patches

State Feature	Description
Jate reature	Menimum height of the notab
Height	Maximum height of the patch
Avg component width	The mean width of the connected components within a patch
Density	Density of foreground pixels within the patch
Aspect ratio	Width/Height of the patch
Gabor filter	8 features capturing the different stroke orientations
Variation of height	Variation in height within a patch
Width variation	Variation in width within a patch
Overlap	Sum of overlap in area between the connected components within a patch
Percentage of text above	Relative location of the patch with respect to the entire document
Number of components	Count of the connected components within a patch
Maximum component size	Maximum size of a component within a patch
Points in convex hull	Number of points in the convex hull of the patch
Maximum run length	The maximum horizontal run length within a patch
Avg run length	The average horizontal run length within a patch
Horizontal Transitions	A count of the number of times the pixel value transitions from white to
	black horizontally
Vertical Transitions	A count of the number of times the pixel value transitions from white to
	black vertically

Features for HW/Print/Noise

Transition Feature	Description
Relative location	Assigned weights based on the relative location - top/bottom or right/left
Convex hull distance	The convex hull distance between the 2 patches
Ratio of aspect ratio	The ratio of the aspect ratio values of the 2 patches
Ratio of number of components	The ratio of the number of components present in the 2 patches

NY English Language Arts Assessment (ELA)-Grade 8

American First Ladies

91 I DE 10

hen George Washington was elected the first President of the United States, there was much discussion about the role his wife, Martha, should play. "Lady Washington," the people began to call her. But Martha did not want to be treated like royalty or aristocracy. She decided that she would be an equal partner with her husband on social occasions. By doing this, she established a very important role for the President's wife-"hostess for the nation."

That role was expanded by James Madison's wife, Dolley Madison, the first woman to be called "First Lady." During the early 1800s, women had few rights. They could not vote or take part in politics. Married women were not even allowed to own property or make a will. A woman's education was usually limited to homemaking skills. But Dolley Madison came from a Quaker family whose community opened its schools to boys and girls. Dolley became an outgoing woman with strong opinions, whose influence on her husband was well known. She was also considered to be the center of society in Washington, D.C., hosting receptions at which politicians and diplomats gathered along with the general public. After her death in 1849, Dolley was honored by President Zachary Taylor, who called her "truly our First Lady for half a century."

Many other First Ladies had strong influences on their husbands, but a woman who really made the job her own was Anna Eleanor Roosevelt, the wife of Franklin Delano Roosevelt. Just two days after her husband's inaugural in 1933, Eleanor held the first press conference ever given by a presidential wife. During FDR's presidency, Eleanor was always there with suggestions, proposals, and ideas. Sharecroppers,1 garment workers, students, and other people whom she had encountered on her travels were brought to Washington to meet the President.



Travel she did-some 38,000 miles in her initial year as First Lady. FDR, weakened by polio, which he had contracted in 1921, was not able to travel easily. Eleanor became the "eyes and ears" of her husband, often making fact-finding trips for him. She saw and inspected everything from prisons to hospitals during those years of the Great Depression.2 She also traveled across the country on lecture tours, wrote articles for magazines, and even wrote a daily newspaper column. During World War II, she became FDR's ambassador to the troops overseas. After her husband's death in 1945, Eleanor was appointed U.S. delegate to the United Nations, where she helped to create the Universal Declaration of Human Rights.

To honor Eleanor Roosevelt's life of service, President Harry Truman called her "The First Lady of the World." At her funeral in 1962, Eleanor Roosevelt's remarkable career was summed up by former presidential candidate Adlai Stevenson: "She would rather light a candle than curse the darkness."

² the Great Depression: a period of tremendous economic hardship in the United States and other countries during the 1930s.

Some Famous First Ladies in History







Martha Washington

(1731 - 1802)

Dolley Madison (1768 - 1849)

Eleanor Roosevelt (1884 - 1962)

Edith Wilson

(1872 - 1961)



Session 2: Part 1

Copyright @ 2001 by New York State Department of Education.

Copyright © 2001 by New York State Department of Education

Session 2: Part 1

Go On

Page 12

Sample Prompt and Answers

How was Martha Washington's role as First Lady different from that of Eleanor Roosevelt? Use information from American First Ladies in your answer.

30 How was Martha Washington's role as First Lady different from that of Eleanor Rooosevelt Use information from "American First Ladies" in your answer.

Martha Washington did not want to be twated like regality on arotocropy. She diade to be an equal parties with her hubband on social occasions. Martia whathisted a very important sole for the Philodents wife - "hosting for the nation".

Cleanor Roosevelt was the First Rady, She became the "lyes and ears" of her husband, often making fact-finding trips for him. She also traveled across the Country on lecture tours, wrote articles for magazines, and even wrote a daily newspaper column. She became FDR'S Anchassador to the troops overseas. Eleanor did nore theigh than Martha Washington.

30 How was Martha Washington's role as First Lady different from that of Eleanor Rooosevu Use information from "American First Ladies" in your answer.

Mathawashington's role as cleared from eleanor first latur Conservent Bearing Martha didn't want to be treated like condite or anstearach. the decided that She wanted to be an aqual partner' with har husban. But Eleana possevent had the tais press conference aver given a prosidential will just two drugs after her husboris in alliguial in 1933. And also Durina 7003 presidency alganos Was allward thank with Dispositions, proposals, and ideas whan her husban food Cacitors the country on lecture tours clarky houspap

30 How was Martha Washington's role as First Lady different from that of Eleanor Rooosevelt Use information from "American First Ladies" in your answer.

One what that martha which affected from Elector Proceed. Because Martha clidat want all that attention clidat what to be treated like regulization Elector did what to be because she traveled Almost Everywhere for Franklin J. B.

> How was Marthe Washington's role as First Lady different from that of Eleanor Roosevelt Dechofmention from American First Ladie's how answer These two is on African Decause Eleanor did some Obort wold man nghts. The seems - Cleanor coold travel for comen gats. But martha did nothing. All she wanted to do is not be provised.

Pag

Page 14

Page 14

Answer Sheet Samples

Martha Washington's Vole was 21 Herent from Elegnor Roosevelt secanse she established a very important role for the president's wire "hostess for the nation" Eleanor Roose velt was different from Marthaberause she was the ambassador of the troops and was decourse plathe United Nations,

Hartha Washington 5 vole as First Lady different was different than Eleanor Busalett'S vole begause Marth washington just went to be parterors with her Misland on Special consister by her husband Social accasions. But Eleanor Proselect stad by her husband The whole way threw.

Relevant LanguageTechnologies all using ML

1. DAR

- Zoning
- Handwriting recognition and interpretation
- 2. NLP and IR
 - 1. Latent Semantic Analysis (LSA)
 - 2. Artificial Neural Network (ANN)
 - 3. Information Extraction (IE)
 - Named entity tagging
 - Profile extraction

Named Entity Tagging

High Scoring Essay

Martha Washington's role as first lady was different from that of Eleanor Roosevelt Martha was called hostess for the nation because she was with her husband during social occasions. But she didn't do anything to help her husband or the U.S. greatly. Eleanor Roosevelt did help a lot. She held the first press conference ever given by a presidential wife. She was always there with suggestions, proposals, and ideas. After Franklin Roosevelt contracted polio, she travelled for him and helped him out. Martha and Eleanor had very different ways of being First Lady.



Output from NE Tagging is XML indicating

- text string
- NE type
- NE subtype

<named entities>

<named entity id="1">Martha Washington

<type>Person</type>

<subtype>Woman</subtype>

</named entity>

<named entities>

<named entity id="2">Eleanor Roosevelt

<type>Person</type>

<subtype>Woman</subtype>

</named entity>

<named entities>

<named entity id="3">Martha

<type>Person</type>

<subtype>Woman</subtype>

</named entity>

cnamed entities>

cnamed entities>

<named entity id="4">U.S>

<type>Location</type> <subtype>Country</subtype>

</named entity>

- Person
- Location
- Verbs triggering events

Summary

- ML is a systematic approach instead of adhockery in designing information processing systems
- Graphical models are useful to deal with complex distributions
- The fully Bayesian approach is computationally expensive but together with Monte Carlo sampling is leading to high performing systems
 - Presented some results in
 - Biometrics/forensics
 - Language-elated technologies