Probabilistic Graphical Models in Machine Learning

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

ICDAR Plenary, Beijing, China September 2011

Plan of Discussion

- Machine Learning (ML)

 History and Problem types solved
- Probabilistic Graphical Models (PGMs)
 - Tutorial
 - Specialized models
- Computational Forensics Application

 Handwriting

What is Machine Learning?

- Automatic construction of programs from examples of input-output behavior
- Marriage of Computer Science and Probability/ Statistics
 - 1. Computer Science:
 - Artificial Intelligence
 - Tasks performed by humans not well described algorithmically
 - Data Explosion
 - User and thing generated
 - 2. Statistics:
 - Methods that learn from data (MLE or Bayesian) 3

When is Machine Learning Needed

- Problems involving uncertainty
 - Perceptual data (images, text, speech, video)
- Information overload
 - Large Volumes
 - Limitations of time, cognitive ability
- Constantly Changing Data Streams
 - Search engine adaptation
- Principled design
 - High performance systems

	ML Problem Types	Example Problems	ML Methods
1	Regression	Keeling Curve, LeToR	Polynomial, GP
2	Classification	OCR, Spam Filter	Logistic Regression, Naïve Bayes
3	Collective Labeling	Speech, HW PoS, NE	HMM, CRF, MEMM
4	Probabilistic Inference	Forensics	BN, MRF
5	Clustering	Data Mining	K-means, EM,

Taxonomy of ML Methods

- First Generation (1960-1980)
 - Perceptrons, Nearest-neighbor, Naïve Bayes
 - Special Hardware, Limited performance
- Second Generation (1980-2000)
 - ANNs, Kalman, HMMs, SVMs
 - HW addresses, speech reco, postal words
 - Difficult to include domain knowledge
 - Black box models fitted to large data sets
- Third Generation (2000-Present)
 - PGMs, Fully Bayesian (including GP)
 - Image segmentation, Text analytics (NE Tagging)
 - Expert prior knowledge with statistical models



USPS-MI OCR

Amperst

USPS-RCR



20 x 20 cell Adaptive Wts

Regression: Learning To Rank Input (x_i) :



46 query-document features Maximum of 124 URLs/query



Output (y): **Relevance Value Target Variable**

_

- Point-wise (0,1,2,3)
- Regression returns continuous value -Allows fine-grained ranking of URLs

Log frequency of query in

of images on page

PageRank of page

URL contains "~"

URL length

Page length

of (out) links on page

Query word in color on page

Traditional IR uses TF/IDF

anchor text

Classification: OCR

Input x={ $x_1...x_{12}$ }: Image Features Output (y): Class Labels { $y^0, y^1, ...y^9$ }















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 1,000 chars/page, 1,000s of pages



Features (x_i) :

Values: Proportion of black pixels in each of 12 cells x_i i=1,..,12

$$\begin{array}{c} x_i^{\ 0} = 0 - 10\% \\ x_i^{\ 1} = 10 - 20\% \\ \dots \end{array} \quad |Val(x_i)| = 10 \end{array}$$

No of parameters=10¹²- 1 Or 1 trillion Per class No of samples needed=??

Probabilistic Inference

Computational Forensics



R = Height Rela-	L = Shape of Loop	A = Shape of	C = Height of	B = Baseline of h	S = Shape of t
tionship of <i>t</i> to <i>h</i>	of <i>h</i>	Arch of h	Cross on t staff		
$r^0 = t$ shorter than h	$l^0 = retraced$	a^0 = rounded	$c^0 = upper half of$	$b^0 = \text{slanting up}$	$s^0 = $ tented
		arch	staff	ward	
$r^1 = t$ even with h	$l^1 = $ curved right side	$a^1 = pointed$	$c^1 = $ lower half of	b^1 = slanting	$s^1 = \text{single stroke}$
	and straight left side		staff	downward	_
$r^2 = t$ taller than h	$l^2 = $ curved left side	a^2 =no set pat-	c^2 = above staff	$b^2 = baseline even$	$s^2 = \text{looped}$
	and straight right side	tern			_
$r^3 = $ no set pattern	l^3 = both sides		$c^3 = \text{no fixed pat-}$	$b^3 = \text{no set pattern}$	$s^3 = closed$
_	curved		tern	_	
	$l^4 =$ no fixed pattern				s^4 = mixture of
	_				shapes

Six Variables, No of parameters= 4,799

Role of PGMs in ML

- Dozens of ML models, Large Data sets
 - PGMs Provide understanding of:
 - model relationships (theory)
 - problem structure (practice)
 - Allow including human knowledge
- Nature of PGMs
 - 1.Represent joint distributions
 - 1.Many variables without full independence
 - 2.Expressive
 - 3.Declarati
 - 2.Inference: Separate model/algorithm errors3.Learning

Representation

BN compression

the the the

R = Height Rela-	L = Shape of Loop	A = Shape of	C = Height of	B = Baseline of h	S = Shape of t
tionship of <i>t</i> to <i>h</i>	of h	Arch of h	Cross on t staff		
$r^0 = t$ shorter than h	$l^0 = retraced$	a^0 = rounded	$c^0 = upper half of$	$b^0 = $ slanting up-	$s^0 = \text{tented}$
		arch	staff	ward	
$r^1 = t$ even with h	$l^1 = curved right side$	$a^1 = pointed$	$c^1 = $ lower half of	b^1 = slanting	$s^1 = \text{single stroke}$
	and straight left side		staff	downward	
$r^2 = t$ taller than h	$l^2 = curved left side$	a^2 =no set pat-	$c^2 = above staff$	$b^2 = baseline even$	$s^2 = \text{looped}$
	and straight right side	tern			
$r^3 = no \text{ set pattern}$	l^3 = both sides		$c^3 = no$ fixed pat-	$b^3 = no \text{ set pattern}$	$s^3 = closed$
_	curved		tern	_	
	$l^4 = no$ fixed pattern				s^4 = mixture of
1	-				shapes

X5

 $x_2 = a_1 x_3 = a$

 $x_2=a,x_3=b$

 $x_2 = a_1 x_3 = c$

 $x_2=b_x_3=a$

 $x_2 = b_1 x_3 = b_2$

 $x_2 = b_1 x_3 = c$

 $x_2 = c_1 x_3 = a$

 $x_2 = c_1 x_3 = b$

 $x_2 = c_1 x_3 = c_2$

 $x_2=d, x_3=a$

 $x_2=d_x_3=b$

 $x_2=d_x_3=c$

 $x_2 = e_1 x_3 = a$

 $x_2 = e_1 x_3 = b$

 $x_2 = e_1 x_3 = c$

а

0.097

0.11

0.09

0.22

0.29

0.14

0.33

0.25

0.20

0.12

0.18

0.17

0.04

0.24

0.09

Six Variables, No of parameters= 4,799

 $\frac{\mathsf{P}(\mathsf{B}|\mathsf{L},\mathsf{A})}{\frac{b}{c}}$

0.59

0.47

0.31

0.44

0.33

0.43

0.17

0.25

0.40

0.41

0.51

0.33

0.61

0.24

0.42

0.22

0.29

0.43

0.22

0.21

0.29

0.33

0.25

0.20

0.35

0.15

0.33

0.25

0.39

0.45

0.091

0.13

0.16

0.11

0.16

0.14

0.17

0.25

0.20

0.12

0.15

0.17

0.11

0.12

0.03



P(L,A,T,B,C,R) = P(L)P(A)P(T) $P(B \mid A,L)P(C \mid T)P(R \mid B,C)$

No of BN parameters=

$$4 + 2 + 59 + 4 + 19 + 63 = 151$$

3% of 4,799

and and and

$\begin{array}{c ccccccccccccccccccccccccccccccccccc$									
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	AN =	AS =	NN =	NS =	NA =	DN =	DS =	DI = Ini-	U =
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	No. of	Forma-	No. of	Forma-	Shape of	No. of	Forma-	tial Stroke	Unusual
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	strokes	tion of	strokes	tion of	Arch of n	strokes	tion of	of d	forma-
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	for for-	Staff of a	for for-	staff of n		for for-	staff of d		tions/use
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $	mation of		mation of			mation of			of symbol
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	a		n			d			
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	$an^0 = one$	$as^0 =$	$nn^0 = one$	$ns^0 =$	$na^0 =$	$dn^0 = one$	$ds^0 =$	$di^0 = top$	$u^0 = \text{for}$
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	continuous	tented	continuous	tented	pointed	continuous	tented	of staff	mation
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	$an^1 = two$	$as^1 = re$ -	$nn^1 =$	$ns^1 = \mathbb{R}^-$	$na^1 =$	$dn^1 = two$	$ds^1 = re$ -	$di^1 = bulb$	$u^1 = sym$ -
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	strokes	traced	two	traced	rounded	strokes	traced		bol
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$			strokes						
$ \begin{array}{ c c c c c c c c c c c c c c c c c c c$	$an^2 =$	$as^2 =$	$nn^2 =$	$ns^2 =$	$na^2 = no$	$dn^2 =$	$ds^2 =$	$di^2 =$	$u^2 = none$
$ \begin{array}{c c c c c c c c c c c c c c c c c c c $	three	looped	three	looped	fixed pat-	three	looped	undeter-	
$ \begin{array}{c ccccccccccccccccccccccccccccccccccc$	strokes	-	strokes	•	tern	strokes	-	mined	
per case staff per case staff per case gle down fixed pattern $an^4 = no$ $as^4 = nn$ $ns^4 = no$ $dn^4 = no$ $ds^4 = sin$	$an^3 = up$ -	$as^3 = no$	$nn^3 = up$ -	$ns^{3} = n0$		$dn^3 = up$ -	$ds^3 = sin$ -	$di^3 = no$	
$an^4 = n0$ $as^4 = nn^4 = n0$ $ns^4 = n0$ $dn^4 = n0$ $ds^4 = sin$	per case	staff	percase	staff		per case	gle down	fixed pat-	
$\boxed{an^4 = no} as^4 = \boxed{nn^4 = no} ns^4 = no} \qquad \boxed{dn^4 = no} ds^4 = sin-$	1						8	tern	
	$an^4 = no$	$as^4 =$	$nn^4 = no$	$ns^4 = no$		$dn^4 = no$	$ds^4 = \sin^4$		
fixed pat- single line fixed pat- fixed pat- fixed pat- gle up	fixed pat-	single line	fixed pat-	fixed pat-		fixed pat-	gle up		
tem down tem tem tem	tern	down	tern	tern		tern			
$as^5 = no$ $ds^5 = no$		$as^5 = no$					$ds^5 = no$		
fixed pat-		fixed pat-					fixed pat-		
tem		tern					tern		

Nine variables. No of parameters = 809,999





Discriminative vs Generative Training

Independent variables $x = \{x_1, .., x_{12}\}$ and binary target y

1. Generative: estimate CPD parameters



2. Discriminative: estimate CRF parameters w_i

Potential Functions (log-linear) $\tilde{P}(y=1 \mid \mathbf{x}) = \exp\left\{\mathbf{w}_0 + \sum_{i=1}^{12} \mathbf{w}_i x_i\right\}$ $\tilde{P}(y=0 \mid \mathbf{x}) = \exp\{0\} = 1$ Unnormalized $\phi_i(x_i, y) = \exp\{w_i x_i | \{y=1\}\},\$ $\phi_0(y) = \exp\{w_0 I\{y=I\}\}\$ $P(y=1 \mid \mathbf{x}) = sigmoid\left\{\mathbf{w}_{0} + \sum_{i=1}^{12} \mathbf{w}_{i} x_{i}\right\} \text{ where } sigmoid(z) = \frac{e^{z}}{1+e^{z}}$ Normalized Naïve I has value *1* Markov Logistic Regression when y=1, else 0 multiclass Jointly optimize 12 parameters $p(y_i | \phi) = y_i(\phi) = \frac{\exp(a_i)}{\sum_i \exp(a_j)}$ High dimensional estimation but correlations accounted for where $a_i = \mathbf{w}_i^{\mathrm{T}} \boldsymbol{\phi}$ Can use much richer features: Edges, image patches sharing same pixels,

Dynamic BN: to model time trajectory

- LeToR relevance values are assigned by human editors
 - Expensive, Can change over time
- Click Logs:
 - Implicit feedback, Current relevance
 - Cheap proxy for editorial label
 - Positional bias: later URLs less likely even if relevant



 C_i : Click on *i*th URL in *retrieved* list

Hidden Variables:

 E_i : did the user *examine* the url? A_i : was the user *attracted* by the url? S_i : was the user *satisfied* by the landing page?

Infer posterior probabilities of hidden vars E_i , A_i and S_i

 $r \equiv P(S_i = 1 \mid E_i = 1) = a_u r_u$

ML estimates of a_u , s_u using EM Forward-backward algorithm to compute posteriors of hidden variables



Collective Labeling: Three Models

Sequence of observations $X = \{X_1, ..., X_k\}$ Need a joint label $Y = \{Y_1, ..., Y_k\}$ Both CRF and MEMM are Discriminative Models Directly obtain P(Y|X)HMM is generative Needs P(X, Y)

Model Trade-offs in expressive power and learnability 1.MEMM and HMM are more easily learned

- Directed models: ML parameter estimates have closed-form
- CRF requires expensive iterative gradient-based approach

2. Ability to use rich feature sets

- HMM needs explicit modeling over features
- CRF and MEMM are discriminative models and avoid this

3.Independence Assumptions made

- MEMM assumes Y_1 independent of X_2 not given Y_2
- Later observation has no effect on current state
 - In <u>activity recognition</u> in video sequence, Frames labelled as running/walking. Earlier frames may be blurry but later ones clearer Model incapable of going back





 $P(Y | X) = \prod_{i=1}^{n} P(Y_i | X_i) P(Y_i | Y_{i-1})$

Inference

Inference and Queries with PGMs

- Inference: Probabilistic Models used to answer queries
- Query Types
 - **1.Probability Queries**
 - Query has two parts
 - *Evidence*: a subset *E* of variables and their instantiation *e*
 - Query Variables: a subset *Y* of random variables in network
 - 2. MAP Queries
 - Maximum a posteriori probability
 - Also called MPE (Most Probable Explanation)

Rare and Common Style Conditional nPRCs

Rare Styles : Looped or tented 't', loop of 'h' with both sides curved



Author:1025c nPRC: 4.17 x 10⁻⁸

Author: 1409c nPRC: 1.23 x 10⁻⁶



Author: 1479b nPRC: 1.66 x 10⁻⁵

Common Styles: Single stroke 't', retraced 'h', pointed arch of 'h', baseline of 'h' slanting down, 't' taller, cross of 't' below

Author: 40b nPRC: 0.76

Author: 130b

Author: 1007c

Author: 685a

PRC = 0.0043 for the dataset. Population size n=100

Rarity Metrics



Complexity of Inference in a BN

Probability Distribution of Evidence

$$P(L,C) = \sum_{A,T,B,R} P(L,A,T,B,C,R)$$
Sum Rule of Probability
$$= \sum_{A,T,B,R} P(L)P(A)P(T)P(B \mid A,L)P(C \mid T)P(R \mid B,C)$$
From the Graphical Model

Probability of Evidence

$$P(L = l^{0}, C = c^{1}) = \sum_{A,T,B,R} P(L = l^{0}) P(A) P(T) P(B \mid A, L) P(C = c^{1} \mid T) P(R \mid B, C = c^{1})$$

More Generally

$$P(E = e) = \sum_{X \setminus E} \prod_{i=1}^{n} P(X_i \mid pa(X_i)) \mid_{E=e}$$

Computed over all *n* variables And all settings of values for them

- An intractable problem
 - #P complete
 - Tractable if tree-width < 25

P: solution in polynomial time NP: verified in polynomial time #P complete: how many solutions

- Approximations are usually sufficient (hence sampling)
 - When P(Y=y|E=e)=0.29292, approximation yields 0.3

Complexity of Inference in a MN

()	(B)-	—(С)	(D	$P(\mathbf{x}) = \frac{1}{Z} \prod_{i=1}^{k} \phi_i(x_p, x_q)$	A 8 0 0 0 0	C D 0 0 0 0 1	φ(A, 1 1	 B) φ(B, 3 3 	$\begin{array}{c c} C & \varphi(C,D) \\ \hline 5 \\ \hline 6 \\ \end{array}$	 <i>P</i> 15 18 22 	P 0.016 0.019	
	А	ssume tł	nat	ead	h variabl	е				where	0 0		1	4	8	32 28	0.033	
	+	kee en	, hir			50	1)			k	0	0 0	2	1	5	10	0.011	
	La	akes only		an	y values:	ξ υ,	13			$7 - \sum_{n=1}^{n} \phi(n - n)$	0 1	0 1	2	1	6	12	0.013	
		Dairy	vico	no	tentials					$\boldsymbol{Z} = \sum \prod \boldsymbol{\varphi}_i(\boldsymbol{x}_p, \boldsymbol{x}_q)$	0 1	1 0	2	9	8	144	0.151	
		Fall W	/150	μo	tentials					x <i>i</i> =1	0 1	. 1 1	2	9	7	126	0.132	
А	в	(A, B)	в	С	$\omega(B,C)$	С	D	ω(C. D)			1 (0 0	3	3	5	45	0.047	
	_	$\Psi(\mathbf{n}, \mathbf{b})$			P (D/D/	_		φ(σ, ε)		×.	1 (0 1	3	3	6	54	0.056	
0	0	1	0	0	3	0	0	5	_		1 (1 0	3	4	8	96	0.100	
											1 (1 1	3	4	7	84	0.088	
0	1	2	0	1	4	0	1	6		/	1 1	0 0	2	1	5	10	0.011	No of Computations
1	0	2	1	0	1	1	0	0		Informa	1 1	0 1	2	1	6	12	0.013	
T	0	5	T	0	1	T	0	0			1 1	1 0	2	9	8	144	0.151	is Exponential in n
1	1	2	1	1	9	1	1	7		Algorithm ^{Z=956}	1 1	. 1 1	2	9	7	126	0.132	

<u>Belief Propagation</u> based on Clique tree (polynomial in *n*)



Based on belief message passing between cliques (nodes) Accurate only when no loops Convergence only for trees

Problem of constructing optimal clique tree

Particle-Based Approximate Inference

Rao-Blackwell: Sampling some variables to perform tractable inference over others

Learning

Learning PGMs

Parameter Learning (given structure)

Bayesian Networks	Markov Networks
Local normalization within each CPD	Global normalization constant (the partition function)
Estimate local groups of parameters separately	Global parameter coupling across the network (even MLE has no closed form)

- Structure Learning
 - Search through network space
- Partial Data

-EM



Parameter Learning For BN



Max Likelihood Est $P(x_5|x_1,x_2)$

	<i>X</i> ₅ = 0	<i>X</i> ₅ = 1	<i>X</i> ₅ = 2	<i>X</i> ₅ = 3
$X_1 = 0, X_2 = 0$	0.50	0	0	0.50
$X_1 = 0, X_2 = 1$	0	1.00	0	0
$X_1 = 0, X_2 = 2$	0.18	0.36	0.27	0.18
$X_1 = 0, X_2 = 3$	0.27	0.40	0.30	0.03
$X_1 = 0, X_2 = 4$	0.22	0.45	0.28	0.05
$X_1 = 1, X_2 = 0$	0.43	0	0.28	0.29
$X_1 = 1, X_2 = 1$	NaN	NaN	NaN	NaN
$X_1 = 1, X_2 = 2$	0.39	0.06	0.33	0.22
$X_1 = 1, X_2 = 3$	0.33	0.17	0.33	0.17
$X_1 = 1, X_2 = 4$	0.42	0.11	0.29	0.18

Bayesian Estimate

	<i>X</i> ₅ = 0			
$X_1 = 0, X_2 = 0$	0.29	0.14	0.29	0.29
$X_1 = 0, X_2 = 1$	0.25	0.25	0.25	0.25
$X_1 = 0, X_2 = 2$	0.25	0.38	0.25	0.12
$X_1 = 0, X_2 = 3$	0.22	0.41	0.31	0.06
$X_1 = 0, X_2 = 4$	0.16	0.52	0.25	0.07
$X_1 = 1, X_2 = 0$	0.29	0.14	0.29	0.29
$X_1 = 1, X_2 = 1$	0.25	0.25	0.25	0.25
$X_1 = 1, X_2 = 2$	0.37	0.05	0.47	0.11
$X_1 = 1, X_2 = 3$	0.33	0.22	0.33	0.11
$X_1 = 1, X_2 = 4$	0.38	0.13	0.29	0.20

Bayesian Estimation

Prior

 $\boldsymbol{\theta} \sim \text{Dirichlet}(\alpha_1, \dots, \alpha_k) \ \alpha_1 = \dots = \alpha_k = 1$

Likelihood $O = \{o_1, ..., o_k\} \sim \text{Multinomial}(\theta_1, ..., \theta_k)$

Posterior $\boldsymbol{\theta}|O \sim \text{Dirichlet}(\alpha'_1, ..., \alpha'_k)$

$$\alpha'_i = \alpha_i + o_i$$
, for $i = 1, ..., k$

25

Parameter Learning for MN



Joint distribution for pairwise MN $p(X) = \frac{1}{Z} \phi_1(X_1, X_2) \cdot \phi_2(X_1, X_5) \cdot \phi_3(X_1, X_6)$ $\cdot \phi_4(X_2, X_5) \cdot \phi_5(X_2, X_6) \cdot \phi_6(X_2, X_3) \cdot \phi_7(X_4, X_6)$ $\cdot \phi_1(X_1) \cdot \phi_2(X_2) \cdot \phi_3(X_3) \cdot \phi_4(X_4) \cdot \phi_5(X_5) \cdot \phi_6(X_6)$

No of Parameters θ_i :

e.g.,

20 + 16 + 20 + 20 + 25 + 15 + 20 + 4 + 5 + 3 + 4 + 4 + 5 = 161

$$\theta_{21} = \log \phi_1(X_1 = 0, X_5 = 0)$$

$$\theta_{22} = \log \phi_1(X_1 = 0, X_5 = 1)$$

$$\vdots$$

$$\theta_{36} = \log \phi_1(X_1 = 3, X_5 = 3)$$



Estimated edge potential for

	Edge Potentia I	$X_{5} = 0$	$X_{5} = I$	X ₅ = 2	$X_5 = 3$	
	$X_I = 0$	3.35	26.10	5.99	1.42	\mathbf{x}
	$X_I = I$	1.54	2.14	1.76	0.94	
	$X_I = 2$	20.01	69.75	33.49	14.90	
	$X_I = 3$	2.99	9.12	4.71	2.25	
og-likelihood	4000 4000 4000 4000 4000 4000 4000 400					

Log-linear model

$$P(x_1,...,x_n:\theta) = \frac{1}{Z(\theta)} \exp\left\{\sum_{i=1}^k \theta_i f_i(D_i)\right\}$$

n: # variables, *k*: # cliques θ_i : parameters

Log-likelihood of *M* i.i.d. samples

$$\ell(\boldsymbol{\theta}) = \sum_{i=1}^{k} \theta_i \left(\sum_{m} f_i(\boldsymbol{\xi}[m]) \right) - M \ln \sum_{\boldsymbol{\xi}} \exp\left(\sum_{i=1}^{k} \theta_i f_i(\boldsymbol{\xi}) \right)$$

Gradient of log-likelihood

$$\frac{\partial}{\partial \theta_i} \frac{1}{M} \ell(\theta) = \frac{\sum_m f_i(\xi[m])}{M} - \frac{\sum_{\xi} f_i(\xi) \exp\left(\sum_{i=1}^k \theta_i f_i(\xi)\right)}{\sum_{\xi} \exp\left(\sum_{i=1}^k \theta_i f_i(\xi)\right)}$$

Concave upward, *BUT* no analytical maximum => Use iterative gradient ascent

Inference step for Z:

Computes unnormalized prob for every *X* setting

=> expensive

x₁

- Approximate inference
 - particle-based methods (MCMC)
 - global algorithm (belief prop)
- Approximate objective
 - Not as much inference
 - Pseudo-likelihood, maxent

Structure Learning of BNs

Compute χ^2 statistics for all pairwise variables & construct En

Let $G = \{V, E\}, E = \{\}, k = 1$

Compute S₀

Structure from

Independence

Chi-squared

- Problem: Many perfect maps for distribution P*
- Goal: Asymptotically recover G*'s equivalence class
- Search through space of BNs
 - Score function for each BN
 - Score₁ (G : D) = log-likelihood (θ_G : D)
 - θ_{G} are parameters of G



Structure Learning of MNs

0.20

0.37

0.04

0.05

0.02

0.08

0.03

0.87

Information-theoretic Chow-Liu algorithm

Algorithm for structure learning:

- 1. Estimate empirical probability: $P(X_1 = N) = \frac{\sum_D \mathbb{1}[X_1 = N]}{\sum_D \mathbb{1}}$
- 2. Calculate all marginal entropies:

$$H(X_1) = -\sum_{X_1} P(X_1) \log(P(X_1))$$

and all pair-joint entropies:

$$H(X_1, X_2) = -\sum_{X_1, X_2} P(X_1, X_2) \log(P(X_1, X_2))$$

- 3. Calculate mutual information: $I(X_1, X_2) = H(X_1) + H(X_2) - H(X_1, X_2) \checkmark$
- 4. Include edges (X_i, X_j) in to the structure if $I(X_1, X_2) \ge threshold$

							I	(X1,.	X_2	$2) \geq t$	hreshold
							G	Sive	S (grap	h
	0.05	0.10	0.00	0.29	0.17				(
	0.00	0.62	0.36	0.24	0.04					1	
	0.12	0.28	0.30	0.30	0.07					\top	
	0.08		0.34	0.17	0.71						\backslash
	0.75				0.00						
						$\left(\right)$	_		(2	
							5		ſ	<u> ۲</u>	
										Ť	Ť
									(3	
	1.34	2.14	2.20	2.4	1 2	2.65	2.1	7	\langle	3	4
	2.14	0.84	1.69	1.92	2 2	2.16	1.6	6			
	2.20	1.69	0.89	1.99		2.23	1.74				
	2.42	1.92	1.99	1.1(о 2	2.44	1.8	9			
	2.66	2.16	2.23	2.44	4	1.36	2.2	!1			
	2.17	1.66	1.75	1.89	9 2	2.20	0.8	7			
			X2							X ₆	
7	X ₁	1.33	0.03		0.02	0.0)2	0.03		0.04	
	X ₂	0.03	0.83	0.5	0.03	0.0	11	0.03		0.05	
	X,	0.02	0.03	0.	0.01	1.1	10	0.02		0.09	
	X ₅	0.03	0.03		0.02	0.0)2	1.36		0.03	28
	-					_					

Summary and Conclusion

- Machine Learning
 - Several generations, with beginnings in DAR field
 - Necessary for changing high volume data
 - To classify, regress, infer, collectively label
- PGMs able to handle complexity
 - BN and MN are expressive
 - Allow incorporating domain knowledge
 - Provide relationships between models
- Computational Forensics Application
 - Handwriting rarity is inferred from PGMs