

Machine Learning using Bayesian Approaches

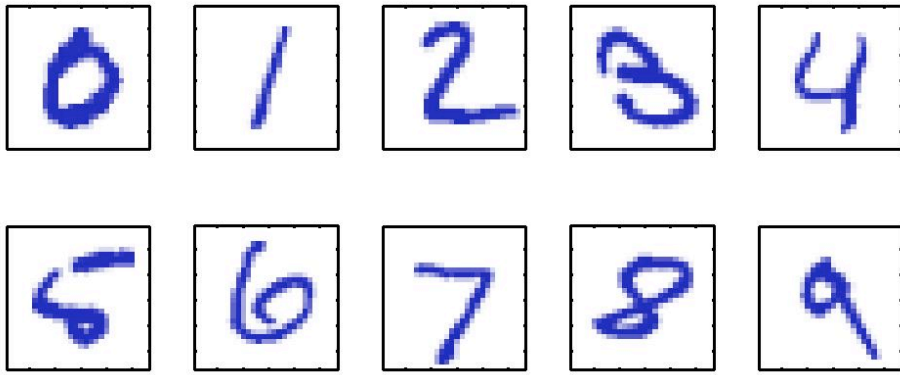
Sargur N. Srihari

University at Buffalo, State University of New York

Outline

1. Progress in ML and PR
2. Fully Bayesian Approach
 1. Probability theory
 - Bayes theorem, Conjugate Distributions
 2. Making predictions using marginalization
 3. Role of sampling methods
 4. Generalization of neural networks, SVM
3. Biometric/Forensic Application
4. Concluding Remarks

Example of Machine Learning: 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
- Provide a probabilistic prediction

Progress in PRML theory

- Unified theory involving
 - Regression and classification
 - Linear models, neural networks, SVMs
 - Basis vectors
 - polynomial, non-linear functions of weighted inputs, centered around samples, etc
- The fully Bayesian approach
- Models from statistical physics
 -

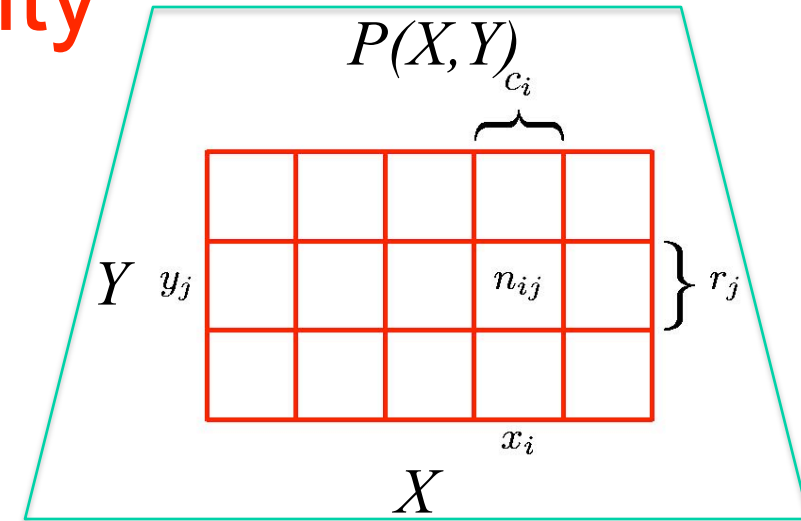
The Fully Bayesian Approach

- 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

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

- 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 \propto likelihood \times prior



Bayesian Probabilities

- Classical or Frequentist view of Probabilities
 - Probability is frequency of random, repeatable event
- Bayesian View
 - Probability is a quantification of uncertainty
 - Whether Shakespeare's plays were written by Francis Bacon
 - Whether moon was once in its own orbit around the sun
- Use of probability to represent uncertainty is not an ad-hoc choice
 - Cox's Theorem:
 - If numerical values represent degrees of belief, then simple set of axioms for manipulating degrees of belief leads to sum and product rules of probability

Role of Likelihood Function

- Uncertainty in w expressed as

$$p(w | D) = \frac{p(D | w)p(w)}{p(D)}$$

where $p(D) = \int p(D | w)p(w)dw$ by Sum Rule

- Denominator is normalization factor
 - Involves marginalization over w
- Bayes theorem in words

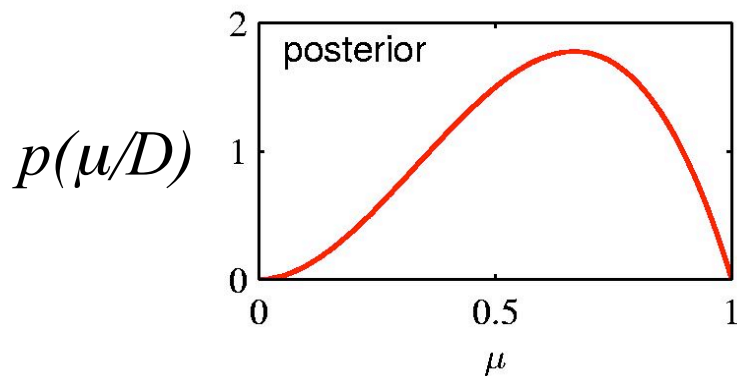
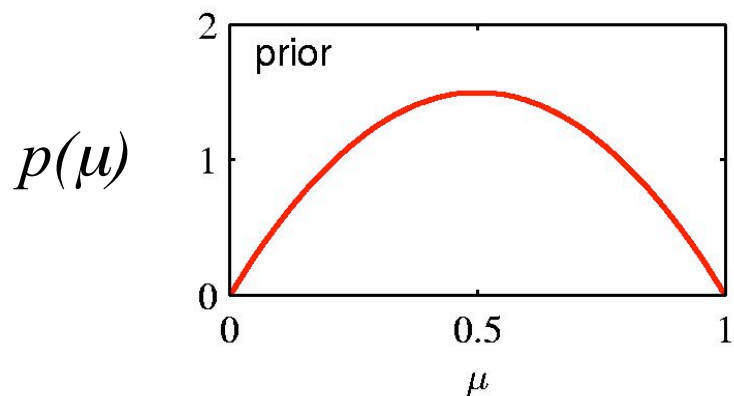
$$posterior \propto likelihood \times prior$$

- Likelihood Function used in both Bayesian and frequentist paradigms
 - Frequentist: w is fixed parameter determined by estimator; error bars on estimate from possible data sets D
 - Bayesian: there is a single data set D , uncertainty expressed as probability distribution over w

Bayesian versus Frequentist Approach

- Inclusion of prior knowledge arises naturally
- Coin Toss Example
 - Suppose we need to learn from coin tosses so that we can predict whether a coin toss will land Heads
 - Fair looking coin is tossed three times and lands Head each time
 - Classical m.l.e of the probability of landing heads is 1 implying all future tosses will land Heads
 - Bayesian approach with reasonable prior will lead to less extreme conclusion

Beta Distribution has property of conjugacy



Bayesian Inference with Beta

- Posterior obtained by multiplying beta prior with binomial likelihood $\binom{N}{m} \mu^m (1-\mu)^{N-m}$ yields

$$p(\mu \mid m, l, a, b) \propto \mu^{m+a-1} (1-\mu)^{l+b-1}$$

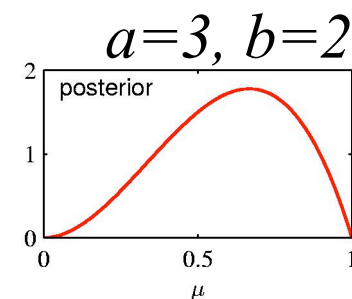
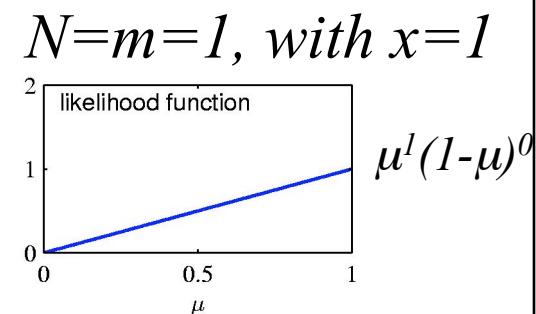
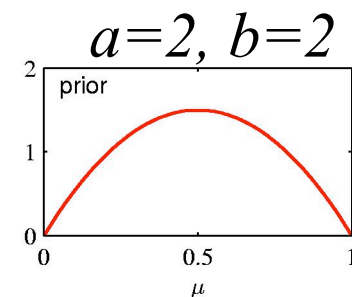
- where $l=N-m$, which is no of tails
- m is no of heads

- It is another beta distribution

$$p(\mu \mid m, l, a, b) = \frac{\Gamma(m+a+l+b)}{\Gamma(m+a)\Gamma(l+b)} \mu^{m+a-1} (1-\mu)^{l+b-1}$$

- Effectively increase value of a by m and b by l
- As number of observations increases distribution becomes more peaked

Illustration of one step in process



Predicting next trial outcome

- Need predictive distribution of x given observed D
 - From sum and products rule

$$\begin{aligned} p(x = 1 | D) &= \int_0^1 p(x = 1, \mu | D) d\mu = \int_0^1 p(x = 1 | \mu) p(\mu | D) d\mu = \\ &= \int_0^1 \mu p(\mu | D) d\mu = E[\mu | D] \end{aligned}$$

- Expected value of the posterior distribution can be shown to be

$$p(x = 1 | D) = \frac{m + a}{m + a + l + b}$$

- Which is fraction of observations (both fictitious and real) that correspond to $x=l$
- Maximum likelihood and Bayesian results agree in the limit of infinite observations
 - On average uncertainty (variance) decreases with observed data

Bayesian Regression (Curve Fitting)

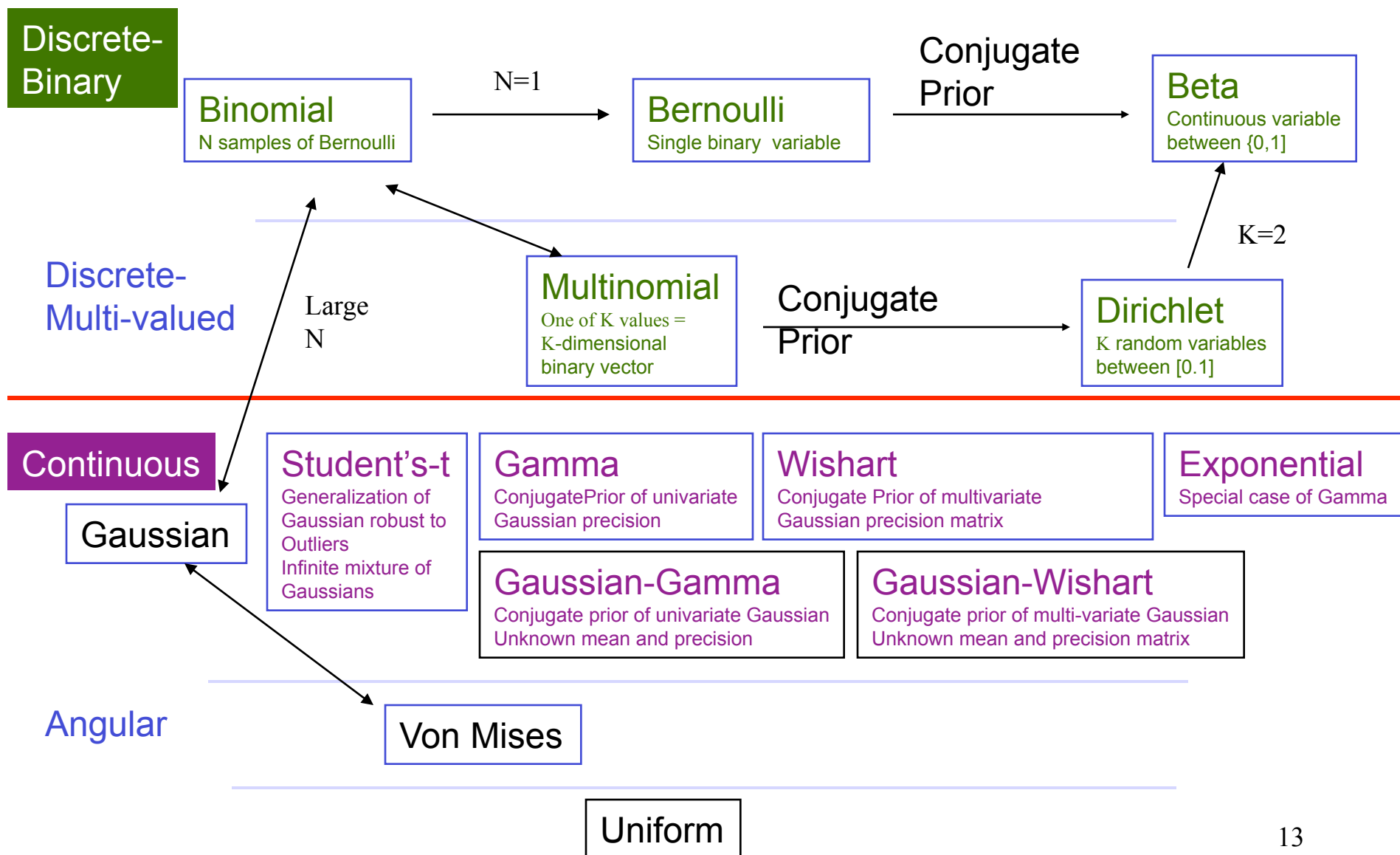
- Estimating polynomial weight vector w
 - In fully Bayesian approach consistently apply sum and product rules and integrate over all values of w
- Given training data x and t and new test point x , goal is to predict value of t
 - i.e, wish to evaluate *predictive distribution* $p(t|x,x,t)$
- Applying sum and product rules of probability

$$\begin{aligned}
 p(t | x, x, t) &= \int p(t, w | x, x, t) dw && \text{by Sum Rule (marginalizing over } w) \\
 &= \int p(t | x, w, x, t) p(w | x, x, t) && \text{by Product Rule} \\
 &= \int \underbrace{p(t | x, w)}_{\text{Posterior distribution over parameters}} \underbrace{p(w | x, t)}_{\text{Also a Gaussian}} dw && \text{by eliminating unnecessary variables}
 \end{aligned}$$

$$p(t | x, w) = N(t | y(x, w), \beta^{-1})$$

Posterior distribution over parameters
Also a Gaussian

Conjugate Distributions



Bayesian Approaches in Classical Pattern Recognition

- Bayesian Neural Networks
 - Classical neural networks focus on maximum likelihood to determine network parameters (weights and biases)
 - Marginalizes over distribution of parameters in order to make prediction
- Bayesian Support Vector Machines

Practicality of Bayesian Approach

- Marginalization over whole parameter space is required to make predictions or compare models
- Factors making it practical:
 - Sampling Methods such as Markov Chain Monte Carlo,
 - Increased speed and memory of computers
- Deterministic approximation schemes
 - Variational Bayes is an alternative to sampling

Markov Chain Monte Carlo (MCMC)

- Rejection sampling and importance sampling for evaluating expectations of functions
 - suffer from severe limitations, particularly with high dimensionality
- MCMC is a very general and powerful framework
 - Markov refers to sequence of samples rather than the model being Markovian
- Allows sampling from large class of distributions
- Scales well with dimensionality
- MCMC origin is in statistical physics (Metropolis 1949)

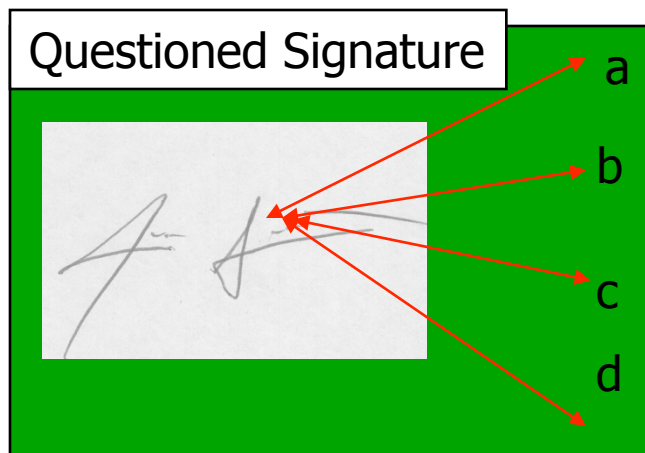
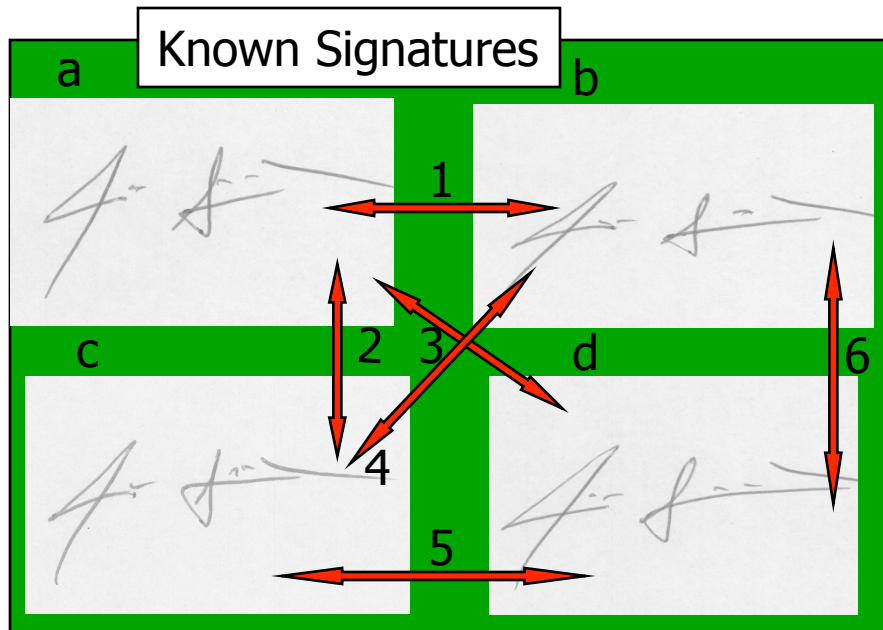
Gibbs Sampling

- A simple and widely applicable Markov Chain Monte Carlo algorithm
- It is a special case of Metropolis-Hastings algorithm
- Consider distribution $p(\mathbf{z}) = p(z_1, \dots, z_M)$ from which we wish to sample
- We have chosen an initial state for the Markov chain
- Each step involves replacing value of one variable by a value drawn from $p(z_i | \mathbf{z}_{\setminus i})$
where symbol $\mathbf{z}_{\setminus i}$ denotes z_1, \dots, z_M with z_i omitted
- Procedure is repeated by cycling through variables in some particular order

Machine Learning in forensics

- Q-K Comparison
 - involves significant uncertainty in some modalities
 - E.g., handwriting
- Learn intra-class variations
 - so as to determine whether questioned sample is genuine

Genuine and Questioned Distributions



Knowns
Compared with
each other

Genuine
Distribution

.4,.1,.2,.3,.2,.1

**Distribution
Comparison**

Questioned
Compared with
Knowns

0,.4,.1,.2
Questioned
Distribution

Methods for Comparing Distributions

- Kolmogorov- Smirnov Test
- Kullback-Leibler Divergence
- Bayesian Formulation

Kolmogorov Smirnov Test

- To determine if two data sets were drawn from same distribution
- Compute a statistic D
 - maximum value of the absolute difference between two c.d.f.s
 - K-S statistic for two cdfs $S_{N1}(x)$ and $S_{N2}(x)$

$$D = \max_{-\infty < x < \infty} |S_{N1}(x) - S_{N2}(x)|$$

- Mapped to a probability P according to

$$P_{KS} = Q_{KS}(\sqrt{N_e} + 0.12 + (0.11/\sqrt{N_e})D)$$

Kullback-Leibler (KL) Divergence

- If unknown distribution $p(\mathbf{x})$ is modeled by approximating distribution $q(\mathbf{x})$
- Average additional information required to specify value of \mathbf{x} as a result of using $q(\mathbf{x})$ instead of $p(\mathbf{x})$

$$\begin{aligned} KL(p \parallel q) &= -\int p(x) \ln q(x) dx - \left(\int p(x) \ln p(x) dx \right) \\ &= -\int p(x) \ln \left\{ \frac{p(x)}{q(x)} \right\} dx \end{aligned}$$

K-L Properties

- Measure of Dissimilarity of $p(x)$ and $q(x)$
- Non-symmetric: $KL(p||q) \neq KL(q||p)$
 - $KL(p||q) \geq 0$ with equality iff $p(x)=q(x)$
 - Proof with convex functions
- Not a metric
- Convert to a probability using

$$D_{KL} = \sum_{b=1}^B P_b \log\left(\frac{P_b}{Q_b}\right)$$

$$P_{KL} = e^{-\zeta D_{KL}}$$

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 $p(D_1 = D_2)$ given that we only observe samples from them
- Method:
 - Determine $p(D_1 = D_2 = \text{a particular distribution } f)$
 - Marginalize over all f (Non-informative prior)
 - If f is Gaussian we obtain a closed-form solution for marginalization over the family of Gaussians F .

Bayesian Comparison of Distributions

$$\begin{aligned}
 P(D_1 = D_2 \mid S_1, S_2) &= \int_F P(D_1 = D_2 = f \mid S_1, S_2) df && \text{Marginalizing over all } f, \text{ a member of family } F \\
 &= \int_F \frac{p(S_1, S_2 \mid D_1 = D_2 = f) p(D_1 = D_2 = f)}{p(S_1, S_2)} df && \text{Bayes rule} \\
 &= \int_F \frac{p(S_1, S_2 \mid D_1 = D_2 = f)}{p(S_1, S_2)} df && \text{Assuming all distributions equally likely or uniform } p(f) \\
 &= \int_F \frac{p(S_1, S_2 \mid D_1 = D_2 = f)}{\int_F p(S_1, S_2 \mid D_1 = D_2 = f) p(D_1 = D_2 = f) df} df && \text{Marginalizing Denominator over all } f \\
 &= \int_F \frac{p(S_1, S_2 \mid D_1 = D_2 = f)}{\int_F p(S_1, S_2 \mid D_1 = D_2 = f) df} df && \text{Due to uniform } p(f) \\
 &= \int_F \frac{p(S_1, S_2 \mid D_1 = D_2 = f)}{\int_F p(S_1 \mid D_1 = f_1) df_1 \int_F p(S_2 \mid D_2 = f_2) df_2} df && \text{Since } S_1 \text{ and } S_2 \text{ are cond independent given } f \\
 &= \frac{Q(S_1 \cup S_2)}{Q(S_1) \times Q(S_2)} \quad \text{where} \quad Q(S) = \int_F p(S \mid f) df
 \end{aligned}$$

Bayesian Comparison of Distributions (F = Gaussian Family)

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

Evaluation of Q for Gaussian family

For Gaussians
closed form for
 Q is obtainable

$$Q(S) = \int_{-\infty}^{\infty} \int_0^{\infty} p(S | \mu, \sigma) d\mu d\sigma$$

$$Q(S) = 2^{-3/2} \times \pi^{(1-n)/2} \times n^{(-1/2)} \times C^{(1-n/2)} \times \Gamma(n/2 - 1)$$

$$C = \sum_{x \in S} x^2 - \frac{(\sum_{x \in S} x)^2}{n}$$

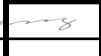





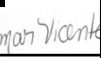
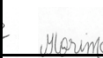

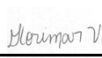





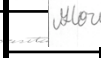
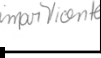
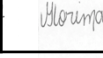
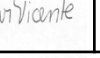
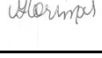
where n is the cardinality of S

Signature Data Set

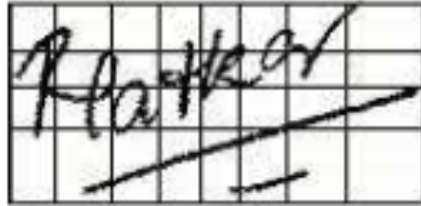
55 individuals

	Horumari Vicente		Q. C. Hall	
			Robert A. Decker	
	Jason J. J.			
				
				
				
				
				
				
				
				
				

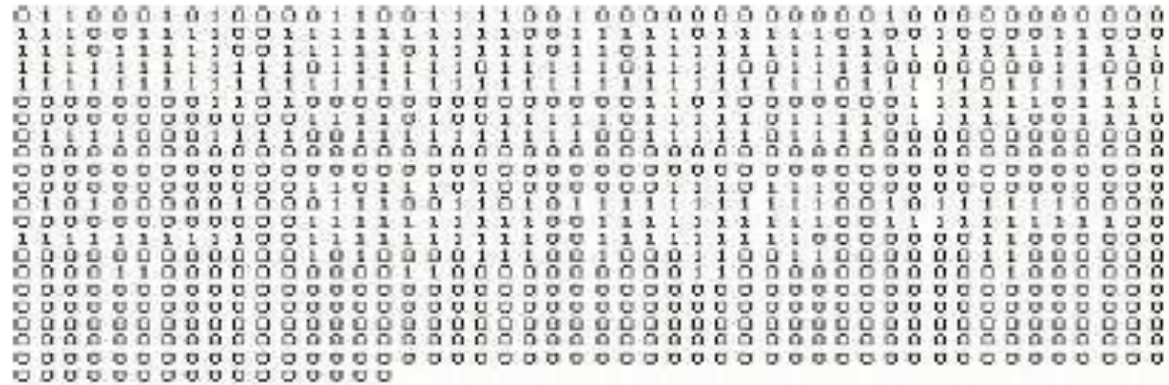
Samples for one person

				
				
<h1>Genuines</h1>				
				
				
<h1>Forgeries</h1>				

Feature Vector and Similarity



(a) variable size grid



(b) binary feature vector

Fig. 2. Feature computation: (a) variable size grid, and (b) binary feature vector.

$$d(X, Y) = \frac{1}{2} \left(1 - \frac{s_{11}s_{00} - s_{10}s_{01}}{[(s_{10} + s_{11})(s_{01} + s_{00})(s_{11} + s_{01})(s_{00} + s_{10})]^{\frac{1}{2}}} \right) \quad (1)$$

where s_{ij} represents the number of corresponding bits of X and Y that have values i and j .

Bayesian Comparison

Clear Separation

Genuines

Forgeries

Signature Verification

Method

☒ Bayesian Comparison ☐ Remove Printed Text

☐ Information-Theoretic

☐ Graph Matching

Select Questioned

Questioned:

Start Testing

Select Writer

1. Author10_1.ws

Author10_1

Knowns

Document ID	Author	Probability of Same
1 Genuine_10_4.png	Author10_1	100.00 %
2 Genuine_10_2.png	Author10_1	99.80 %
3 Genuine_10_5.png	Author10_1	99.50 %
4 Genuine_10_6.png	Author10_1	99.50 %
5 Genuine_10_7.png	Author10_1	99.04 %
6 Genuine_10_3.png	Author10_1	97.59 %
7 Genuine_10_8.png	Author10_1	96.34 %
8 Genuine_10_1.png	Author10_1	93.54 %
9 Forgery_10_2.png	Author10_1	0.62 %
10 Forgery_10_1.png	Author10_1	0.30 %
11 Forgery_10_4.png	Author10_1	0.17 %
12 Forgery_10_8.png	Author10_1	0.16 %
13 Forgery_10_6.png	Author10_1	0.11 %
14 Forgery_10_3.png	Author10_1	0.11 %
15 Forgery_10_7.png	Author10_1	0.05 %
16 Forgery_10_5.png	Author10_1	0.05 %

Tested 16 out of 16

Page up

Page down

Information-Theoretic (KS+KL)

Signature Verification

Method


☐ Bayesian Comparison ☐ Remove Printed Text

☒ Information-Theoretic

☐ Graph Matching

Select Questioned

Questioned:



Start Testing

Select Writer

1. Author10_1.ws


Author10_1

Knowns

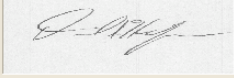
Document ID	Author	Probability of Same
1 Genuine_10_1.png	Author10_1	98.44 %
2 Genuine_10_6.png	Author10_1	97.37 %
3 Genuine_10_2.png	Author10_1	97.37 %
4 Genuine_10_7.png	Author10_1	96.07 %
5 Genuine_10_4.png	Author10_1	96.07 %
6 Genuine_10_8.png	Author10_1	88.32 %
7 Genuine_10_3.png	Author10_1	88.30 %
8 Genuine_10_5.png	Author10_1	88.30 %
9 Forgery_10_2.png	Author10_1	14.87 %
10 Forgery_10_1.png	Author10_1	14.87 %
11 Forgery_10_7.png	Author10_1	4.86 %
12 Forgery_10_6.png	Author10_1	4.86 %
13 Forgery_10_5.png	Author10_1	4.86 %
14 Forgery_10_4.png	Author10_1	4.86 %
15 Forgery_10_3.png	Author10_1	4.86 %
16 Forgery_10_8.png	Author10_1	4.86 %

Tested 16 out of 16

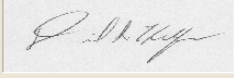
Page up Page down




Genuine_10_1.png



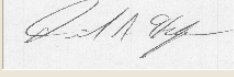
Genuine_10_2.png



Genuine_10_3.png



Genuine_10_4.png



Genuine_10_5.png

Graph Matching

Signature Verification

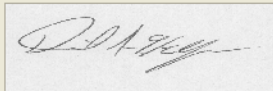
Method:

- ☐ Bayesian Comparison
- ☐ Information-Theoretic
- ☒ Graph Matching

☐ Remove Printed Text

Select Questioned

Questioned:



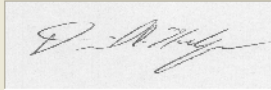
Start Testing

Select Writer

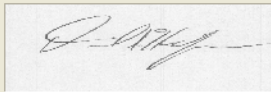
1. Author10_1.ws

Author10_1

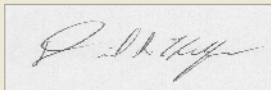
Knowns:




Genuine_10_1.png



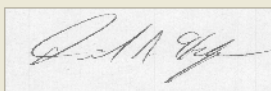
Genuine_10_2.png



Genuine_10_3.png



Genuine_10_4.png



Genuine_10_5.png

Document ID	Author	Probability of Same
1 Genuine_10_5.png	Author10_1	99.96 %
2 Forgery_10_2.png	Author10_1	98.89 %
3 Genuine_10_1.png	Author10_1	98.44 %
4 Forgery_10_1.png	Author10_1	97.05 %
5 Genuine_10_4.png	Author10_1	94.55 %
6 Forgery_10_7.png	Author10_1	91.57 %
7 Genuine_10_8.png	Author10_1	91.12 %
8 Genuine_10_6.png	Author10_1	83.80 %
9 Forgery_10_5.png	Author10_1	78.48 %
10 Forgery_10_6.png	Author10_1	68.07 %
11 Genuine_10_7.png	Author10_1	60.30 %
12 Forgery_10_3.png	Author10_1	56.57 %
13 Genuine_10_3.png	Author10_1	52.39 %
14 Forgery_10_4.png	Author10_1	32.76 %
15 Forgery_10_8.png	Author10_1	14.77 %
16 Genuine_10_2.png	Author10_1	14.52 %

Tested 16 out of 16

Page up

Page down

Comparison of KS and KL

* Threshold Methods

Method	Accuracy
16 Knowns Used	
Distance Threshold	78.50%
Distance Statistics	77.76%
Naïve Bayes	75.00%
One-Class SVM	54.00%

**Information-Theoretic

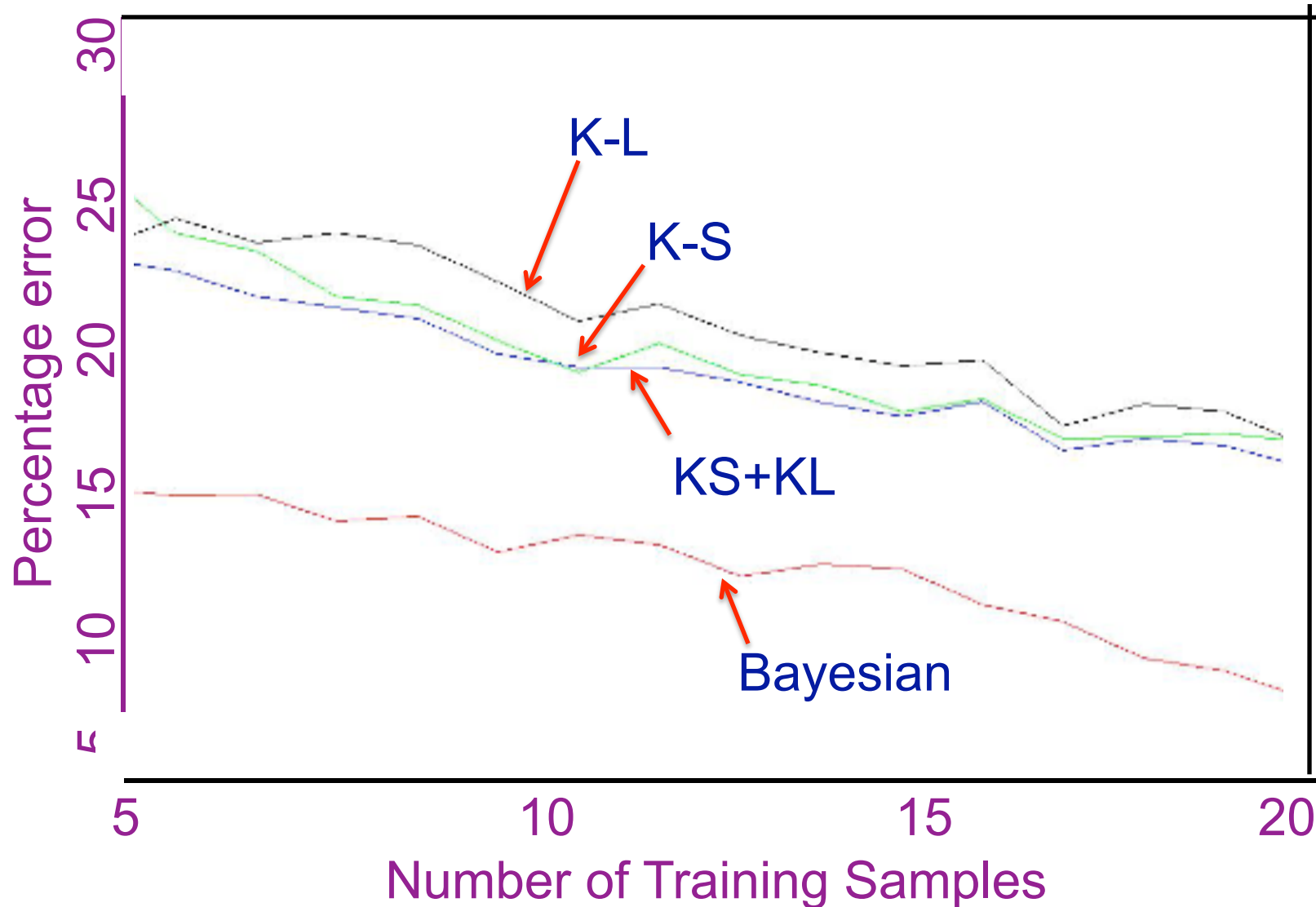
Knowns	KS	KL	KS & KL
4	76.38%	75.10%	79.50%
7	80.10%	79.50%	81.0%
10	83.53%	84.27%	86.78%

*S.N.Srihari A.Xu M.K.Kalera, Learning strategies and classification methods for off-line signature verification.
Proc. of the 7th Int. Workshop on Frontiers in handwriting recognition (IWHR), pages 161-166, 2004

**Harish Srinivasan, Sargur.N.Srihari and Matthew.J.Beal, Machine learning for person identification and verification,
Machine Learning in Document Analysis and Recognition (S. Marinai and H. Fujisawa, ed.), Springer 2008

Error Rates of Three Methods

CEDAR Data Set



Why Does Bayesian Formalism do much better?

- Few samples available for determining distributions
- Significant uncertainty exists
- Bayesian models uncertainty better than information-theoretic method which assumes that distributions are fixed

Summary

- ML is a systematic approach instead of ad-hoc methods in designing information processing systems
- There is much progress in ML and PR
 - developing a unified theory of different machine learning approaches, e.g., neural networks, SVM
- Fully Bayesian approach involves making predictions involving complex distributions, made feasible by sampling methods
- Bayesian approach leads to higher performing systems with small sample sizes

Lecture Slides

- PRML lecture slides at

<http://www.cedar.buffalo.edu/~srihari/CSE574/index.html>