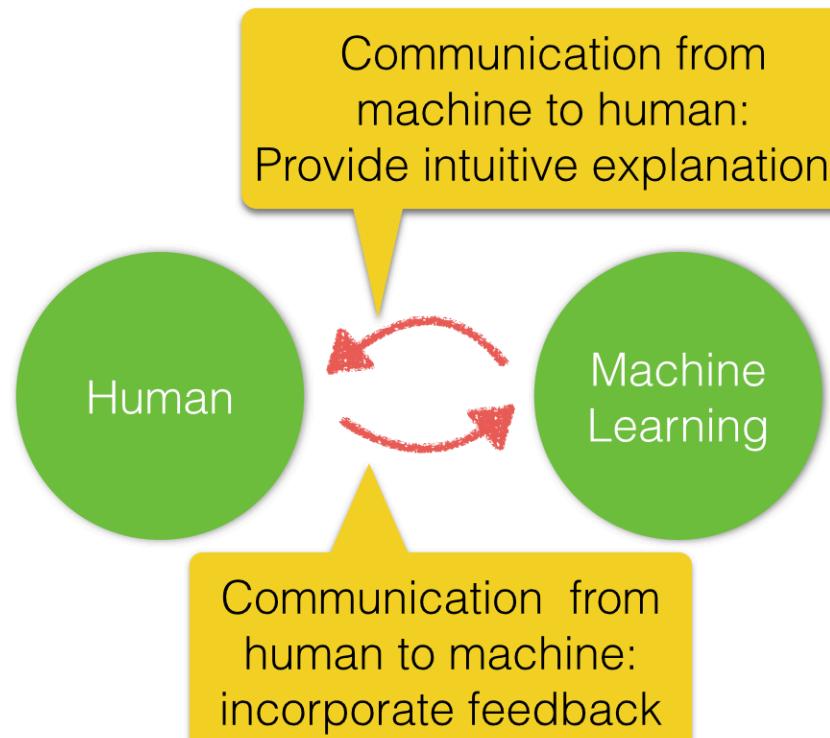


Motivation

Develop interactive machine learning framework to bring together data, machine learning algorithms and human experts' knowledge.



First step: communication from machine to human

Why Case-based Reasoning?

1. Case-based Reasoning (CBR) is intuitive.

Example of CBR (simplified)

- a. My child is sick. He is having fever, nausea and coughs.
- b. Find the illness that has the most similar symptoms → Flu
- c. Apply the solution for flu: rest, ibuprofen, liquid.

2. Exemplar-based reasoning is the way humans think – machine can better support peoples' decision-making by representing data in the same way.

Human's tactical decision is based on exemplar-based reasoning.^[1]

Skilled fire fighters use recognition-primed decision making – a situation is matched to typical cases.^[2,3]



Limitations of CBR

- Does not leverage data
- Requires previous solutions
- Does not scale to complex problems

[1] M.S. Cohen, J.T. Freeman, and S. Wolf. Metarecognition in time-stressed decision making: Recognizing, critiquing, and correcting. *Human Factors*, 1996.

[2] A. Newell and H.A. Simon. Human problem solving. Prentice-Hall Englewood Cliffs, 1972

[3] G.A. Klein. Do decision biases explain too much. *HFES*, 1989.

Bayesian Case Model (BCM): Generative Approach for Case-based Reasoning and Prototype Classification

Been Kim, Cynthia Rudin and Julie Shah
{beenkim, rudin, julie_a_shah}@csail.mit.edu

Bayesian Case Model (BCM)

Leverage the power of examples (**prototypes**) and hot features (**subspaces**) to explain machine learning results

Bayesian generative models

+ Case-based reasoning

= Bayesian Case Model

Prototype

Subspace

Quintessential observation that best represents the cluster

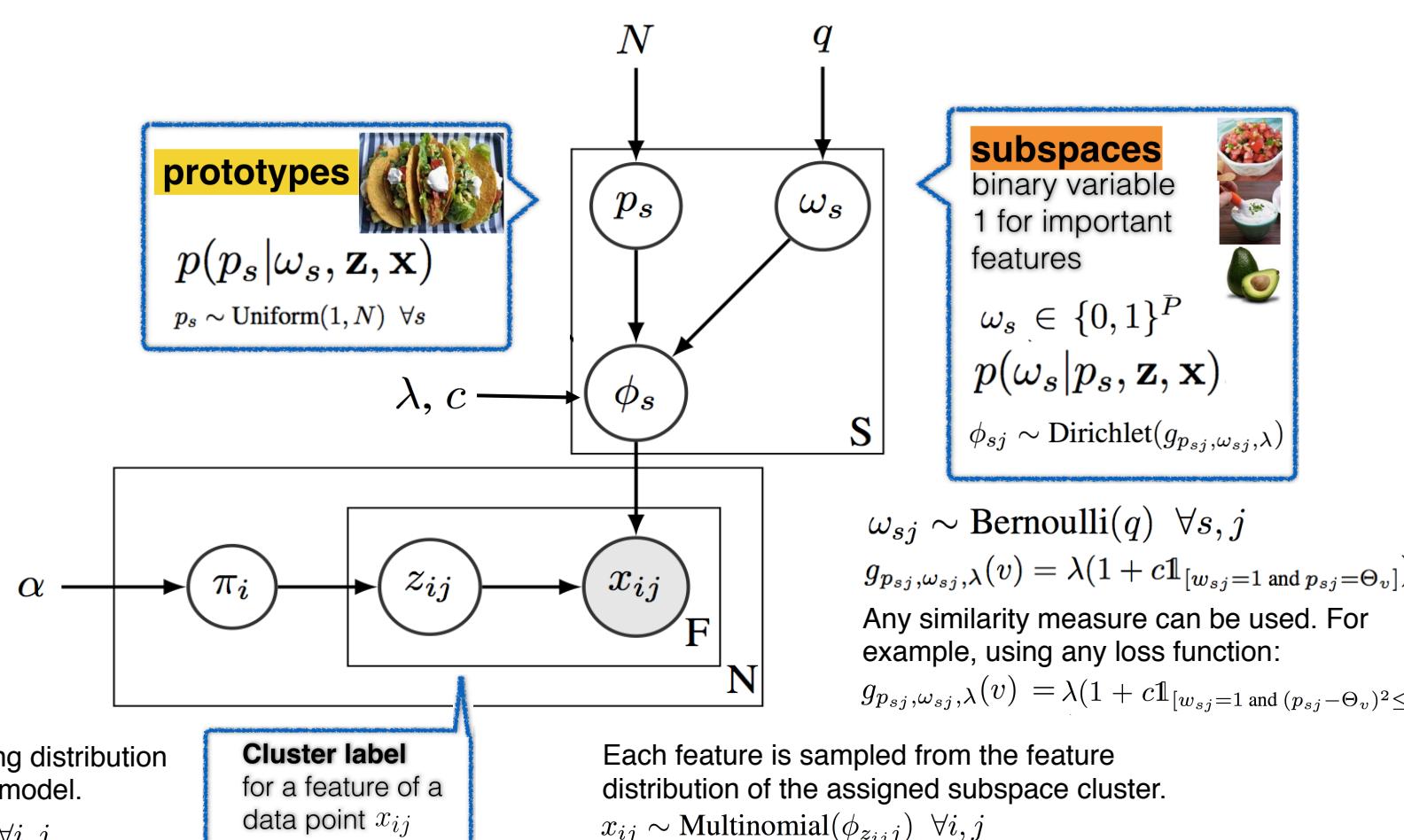
Sets of important features in characterizing clusters



prototypes
 $p(p_s | \omega_s, z, x)$
 $p_s \sim \text{Uniform}(1, N) \quad \forall s$

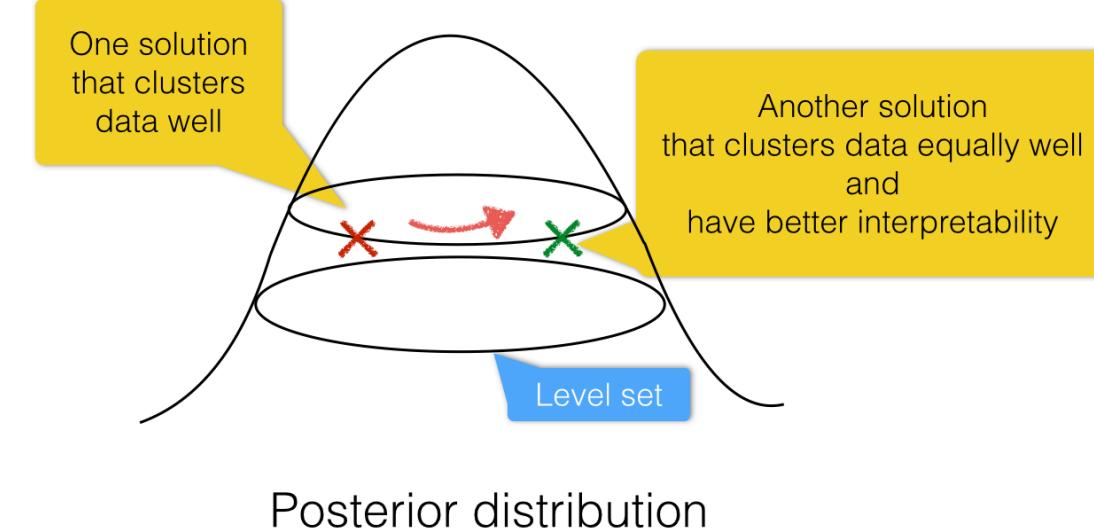
subspaces
binary variable 1 for important features
 $\omega_s \in \{0, 1\}^P$
 $p(\omega_s | p_s, z, x)$
 $\phi_{sj} \sim \text{Dirichlet}(g_{psj}, \omega_{sj}, \lambda)$

subspaces
binary variable 1 for important features
 $\omega_s \in \{0, 1\}^P$
 $p(\omega_s | p_s, z, x)$
 $\phi_{sj} \sim \text{Dirichlet}(g_{psj}, \omega_{sj}, \lambda)$



Improving interpretability without sacrificing performance

Joint inference on prototypes, subspaces and cluster labels is the key.



An example of BCM

prototypes
subspaces

Raw data



Collapsed Gibbs sampling for inference

Integrating out and for efficient inference

$$p(z_{ij} = s | z_{i-j}, x, p, \omega, \alpha, \lambda) \propto \frac{\alpha/S + n_{(s,i,-j,\cdot)}}{\alpha + n} \times \frac{g(p_{sj}, \omega_{sj}, \lambda) + n_{(s,\cdot,j,x_{ij})}}{\sum_s g(p_{sj}, \omega_{sj}, \lambda) + n_{(s,\cdot,j,\cdot)}}$$

where $n_{(s,i,j,v)} = \mathbb{1}(z_{ij} = s, x_{ij} = v)$

$$p(\omega_{sj} = b | q, p_{sj}, \lambda, \phi, x, z, \alpha) \propto \begin{cases} q \times \frac{B(g(p_{sj}, 1, \lambda) + n_{(s,\cdot,j,\cdot)})}{B(g(p_{sj}, 1, \lambda))} & b=1 \\ 1 - q \times \frac{B(g(p_{sj}, 0, \lambda) + n_{(s,\cdot,j,\cdot)})}{B(g(p_{sj}, 0, \lambda))} & b=0 \end{cases}$$

where B is the Beta function



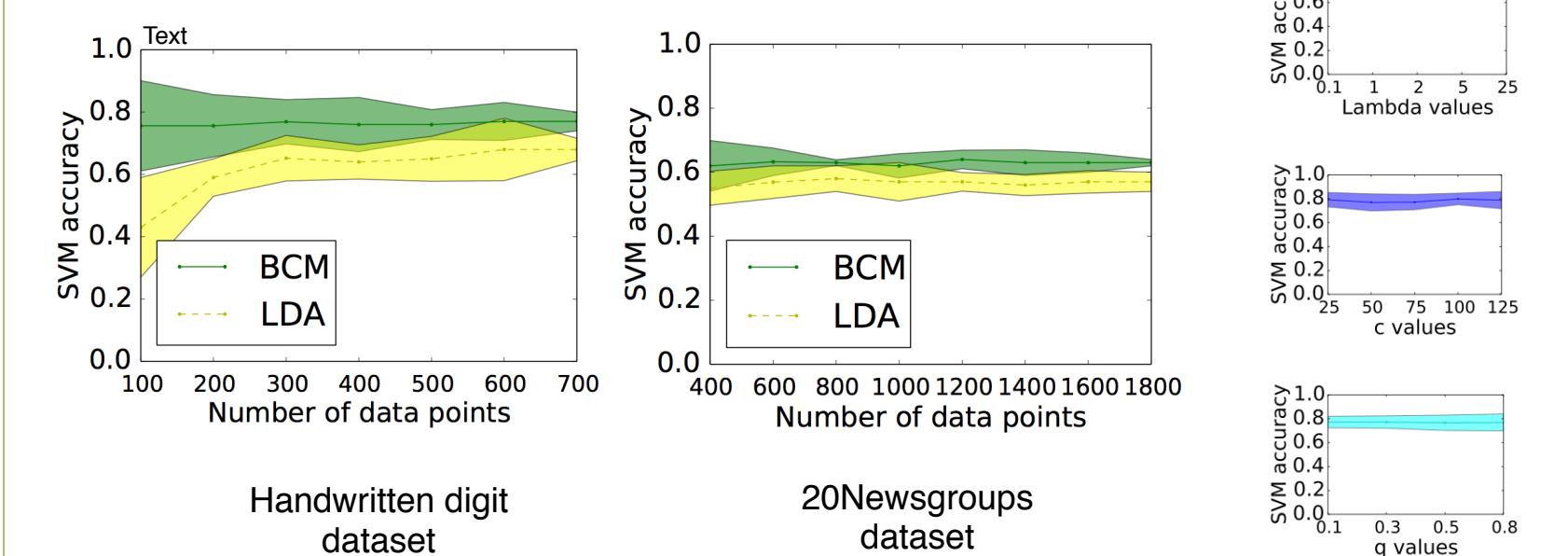
Results

1. Learned prototypes and subspaces

Dataset	Prototypes and Subspaces from BCM	
a. Recipes	Prototypes from BCM (Recipe names)	Subspaces from BCM and Ingredients (subspaces)
Herbs and tomato in pasta	basil, garlic, Italian seasoning, oil, pasta, pepper, salt, tomato	
Generic chili recipe	beer, chili powder, cumin, garlic, meat, oil, onion, pepper, salt, tomato	
Microwave brownies	baking powder, sugar, butter, chocolate, chopped pecans, eggs, flour, salt, vanilla	
Spiced-punch	cinnamon stick, lemon juice, orange juice, pineapple juice, sugar, water, whole cloves	
b. USPS handwritten digits	Prototype A	Subspace A
	Subspace B	Subspace B
	Subspace C	Subspace C
	Subspace D	Subspace D
	Subspace E	Subspace E
	Iteration	

2. BCM maintains accuracy while achieving interpretability

Evaluated using features learned by BCM or LDA + SVM



3. BCM improves humans understanding

- Objective measure of human understanding
- Participant's task is to assign the ingredients of a specific dish (a new data point) to a cluster.
- Statistically significantly better performance with BCM **85.9%** v.s. **71.3%** $p \ll 0.001$

