

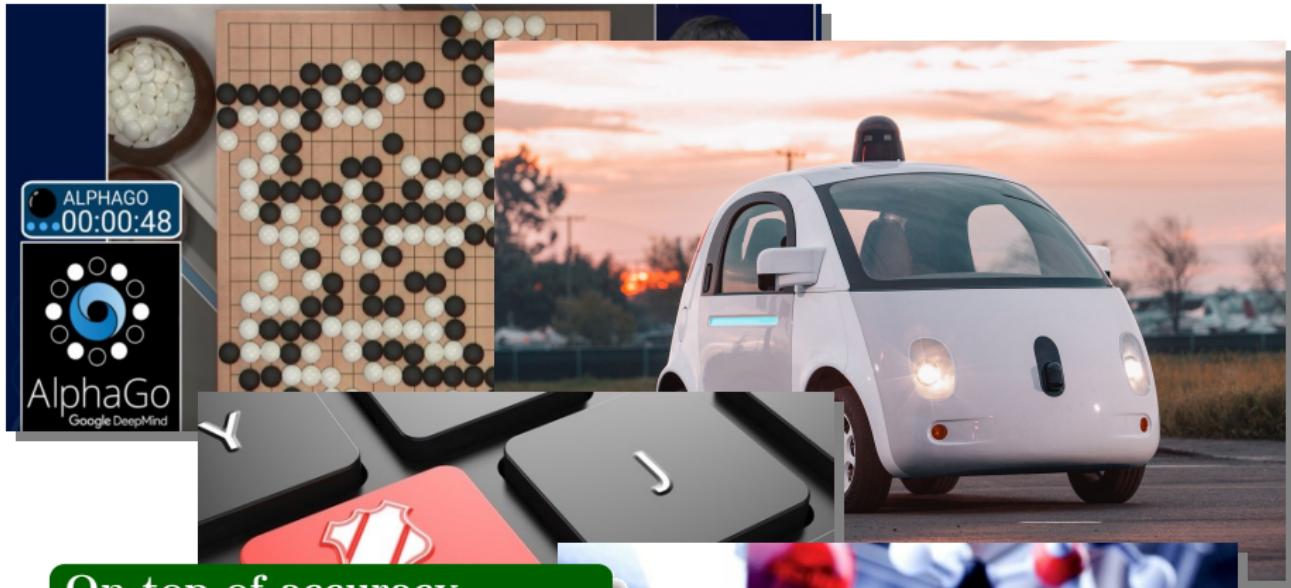
# A Bayesian Non-parametric Approach to Understand World Economies

Melanie F. Pradier, Viktor Stojkoski, Zoran Utkovski,  
Ljupco Kocarev, and Fernando Perez-Cruz

Friday 17, 2017



- High-dimensional count data.
- Focus on **Data Exploration**.



## On top of accuracy...

- Safety
- Fairness
- Understanding

We need **Interpretability**

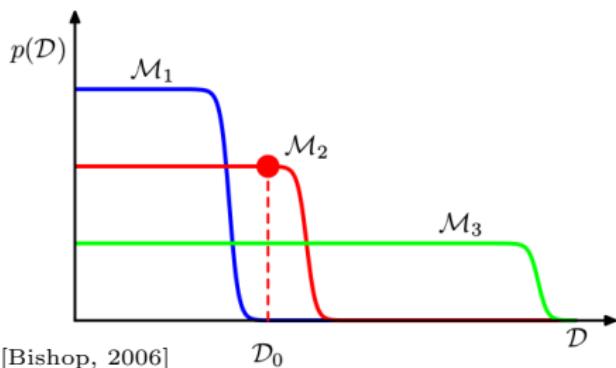
# Interpretability

[F. Doshi-Velez, B. Kim, *Towards A Rigorous Science of Interpretable Machine Learning*]

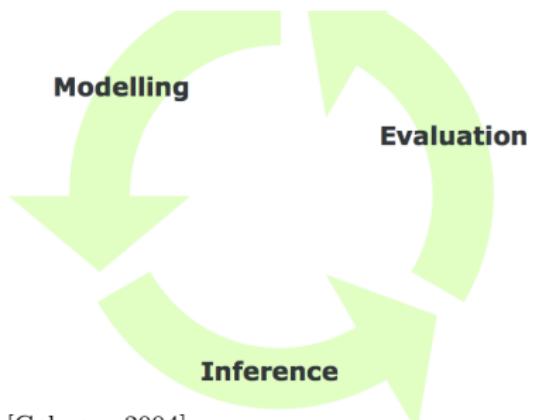
- Understandable for humans (Doshi-Velez, 2017)
- “Right to explanation” (European Legislation, 2018)

## Focus: Data Exploration

- Adequate model assumptions
- Constrained solution spaces



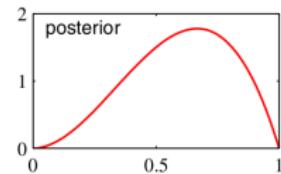
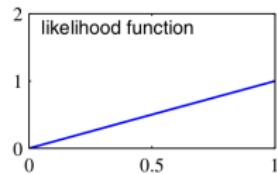
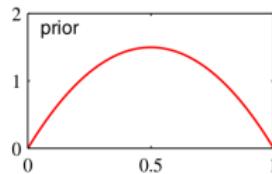
Multidisciplinary Research



[Gelman, 2004]

# Bayesian Non-Parametric Models

- Bayesian: Combine Prior Knowledge with Data Evidence



[Bishop, 2006]

- Non-parametric
  - actually... really large parametric model
  - number of latent variables grows with data

During my PhD...

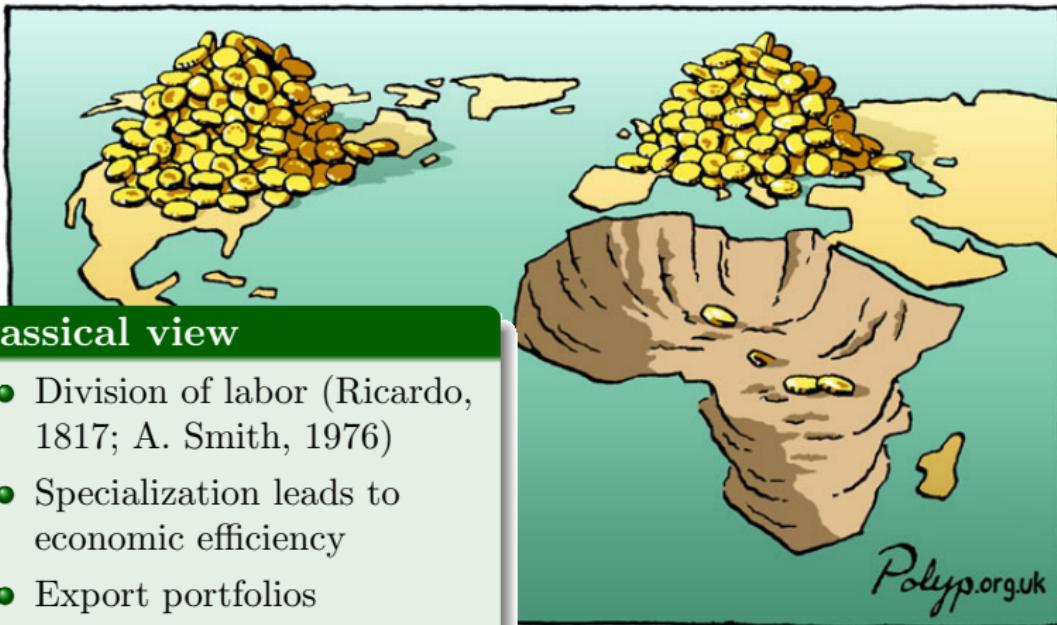
- Marathon modeling
- Biomarker discovery
- **Economic complexity**

# Outline

- ❶ Motivation
- ❷ Theoretical Background
- ❸ Our approach
- ❹ Results

# Motivation: Wealth of Nations

What makes some countries wealthier than others?



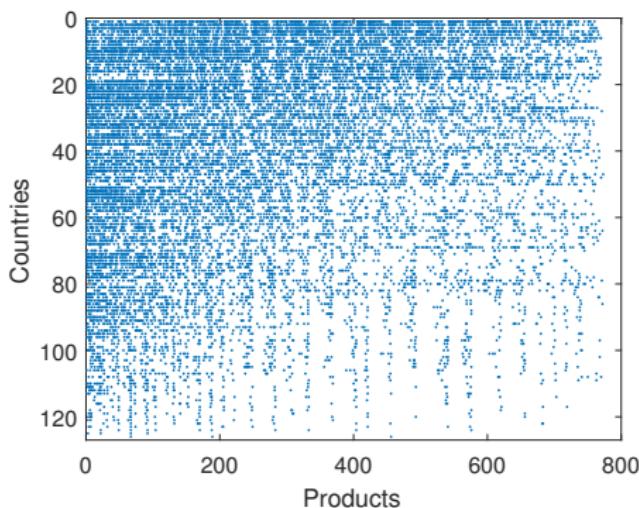
## Classical view

- Division of labor (Ricardo, 1817; A. Smith, 1976)
- Specialization leads to economic efficiency
- Export portfolios

→ block-structure

# Motivation: Wealth of Nations

The reality:



Properties:

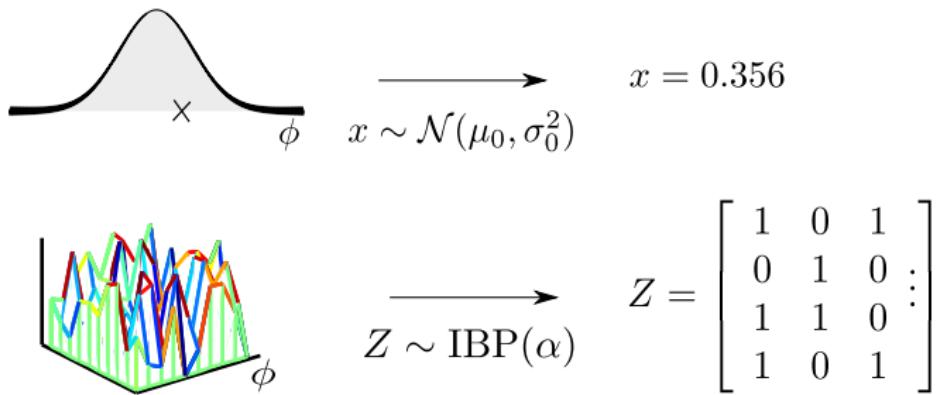
- ① Triangularity
- ②  $D \gg N$

## Our Contribution

Develop an Infinite Poisson-Gamma Model

- Flexible prior
- Feature sparsity

# Indian Buffet Process (Ghahramani et.al, 2006)



- IBP: distribution over binary matrices  $Z_{N \times K}$
- Model chooses number of hidden features,  $K \rightarrow \infty$

# Indian Buffet Process

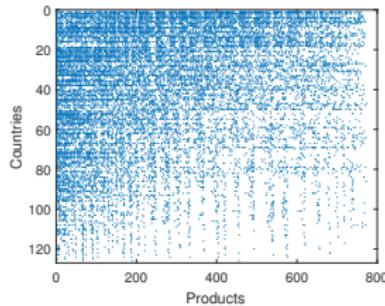
(Slides from F.J.R. Ruiz)

							...
	1	1	0	1	0	1	
	1	0	1	0	0	1	
	0	0	1	0	1	1	
	⋮						

# Our Approach



$$\text{N countries} \quad N \times D \quad D \text{ products} \\ X = p_x \left( \begin{array}{c} N \times K \\ Z \\ \cdot \\ K \times D \\ B \end{array} \right) \quad K \text{ latent features}$$



## Generative Model

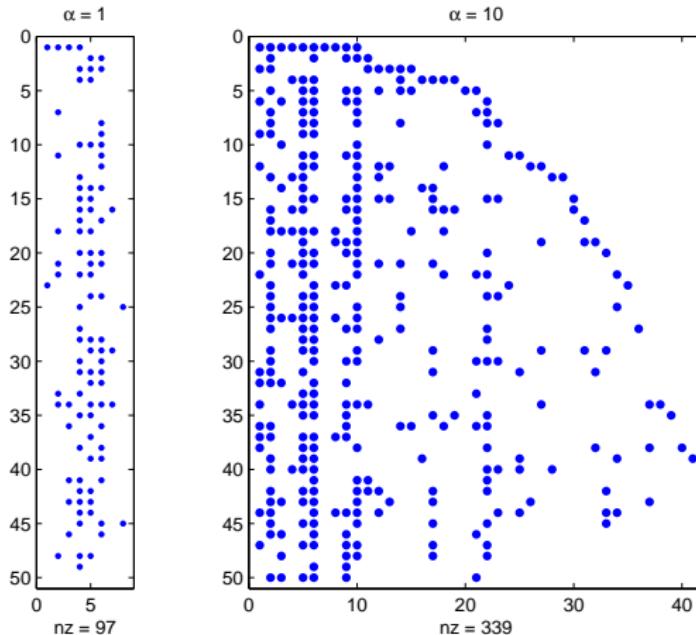
$$x_{nd} \sim \text{Poisson}(\mathbf{Z}_{n\bullet} \mathbf{B}_{\bullet d}) \quad (1)$$

$$B_{kd} \sim \text{Gamma}(\alpha_B, \frac{\mu_B}{\alpha_B}) \quad (2)$$

$$\mathbf{Z}_{n\bullet} \sim \text{IBP}(\alpha) \quad (3)$$

# A limitation of the IBP

- **Disadvantage:** Mass parameter  $\alpha$  couples both  $J_n$  and  $K^+$



# Beyond the standard IBP

## Three-parameter IBP (Teh et.al, 2007)

- More flexible distribution for feature weights

$$\mathbf{Z}_{n\bullet} \sim \text{BeP}(\mu) \quad (4)$$

$$\mu \sim \text{SBP}(1, \alpha, H, \textcolor{red}{c}, \sigma) \quad (5)$$

$$p(J_{new}) \sim \text{Poisson} \left( \alpha \frac{\Gamma(1 + \textcolor{red}{c})\Gamma(n + \textcolor{red}{c} + \sigma - 1)}{\Gamma(n + \textcolor{red}{c})\Gamma(\textcolor{red}{c} + \sigma)} \right)$$



## Restricted IBP (Doshi-Velez et.al, 2015)

- Arbitrary prior  $f$  over  $J_n$

$$\mathbf{Z}_{n\bullet} \sim \text{R-BeP}(\mu, \textcolor{red}{f}) \quad (6)$$

$$\mu \sim \text{BP}(1, \alpha, H) \quad (7)$$

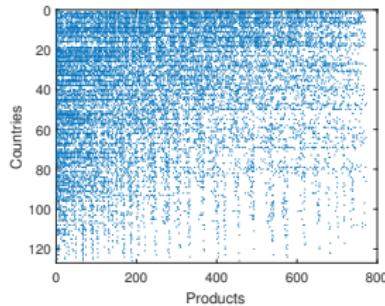


- Combination of both
- Flexible prior

# Our Approach



$$\begin{matrix} N \times D \\ D \text{ products} \\ N \text{ countries} \end{matrix} = p_x \left( \begin{matrix} N \times K \\ N \times K \\ Z \end{matrix} \cdot \begin{matrix} K \times D \\ K \times D \\ B \end{matrix} \right)$$



## Generative Model

$$x_{nd} \sim \text{Poisson}(\mathbf{Z}_{n\bullet} \mathbf{B}_{\bullet d}) \quad (8)$$

$$B_{kd} \sim \text{Gamma}(\alpha_B, \frac{\mu_B}{\alpha_B}) \quad (9)$$

$$\mathbf{Z}_{n\bullet} \sim \text{3R-IBP}(\alpha, c, \sigma, f) \quad (10)$$

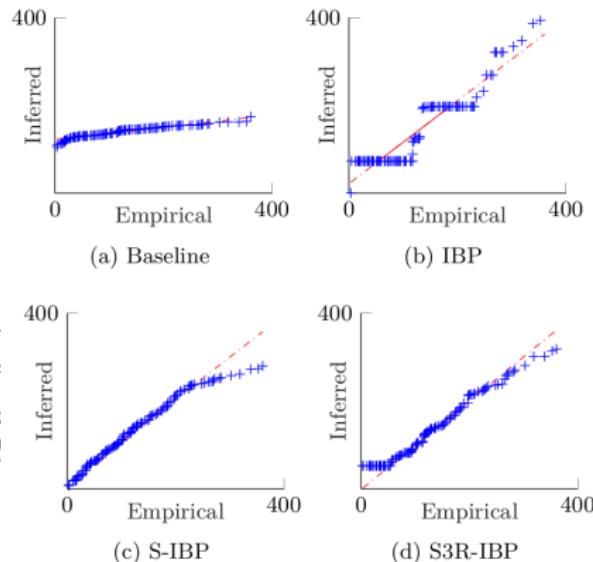
# Results

Metric	IBP	S-IBP	S3R-IBP
Log Perplexity	<b>2.91 ± 0.17</b>	15.29 ± 1.84	3.03 ± 0.04
Coherence	-149.36 ± 7.56	-178.44 ± 4.50	<b>-140.51 ± 2.73</b>
Lift	120.36 ± 3.66	121.77 ± 0.71	<b>125.28 ± 0.50</b>

(a) SITC database ( $N = 126$ ,  $D = 744$ )

Metric	IBP	S-IBP	S3R-IBP
Log Perplexity	<b>2.85 ± 0.06</b>	8.75 ± 0.95	2.88 ± 0.07
Coherence	-148.91 ± 10.57	-168.39 ± 13.16	<b>-134.51 ± 4.4</b>
Lift	152.12 ± 6.85	157.38 ± 2.20	<b>161.80 ± 0.50</b>

(b) HS database ( $N = 123$ ,  $D = 4890$ )



# Results

## Interpretability

F0: Bias	F1: Agriculture	F2: Clothing I	F3: Farming	F4: Clothing II	
Non-Coniferous Worked Wood Bran and Other Cereals Residues Misc. Non-Iron Waste	Vegetables Fruit or Vegetable Juices Misc. Fruit	Synthetic Knitted Undergarments Misc. Feminine Outerwear Misc. Knitted Outerwear	Misc. Animal Oils Bovine and Equine Entrails Bovine meat	Synthetic Woven Fabrics Non-retail Synthetic Yarn Woven Fabric < 85% Discontinuous Synthetic Fibres	
F5: Electronics I	F6: Processed Materials	F7: Electronics II	F8: Materials I	F9: Machinery I	
Misc. Electrical Machinery Vehicles Stereos Misc. Data Processing Equipment	Baked Goods Metal Containers Misc. Edibles	Measuring Controlling Instruments Mathematical Calculation Instruments Misc. Electrical Instruments	Misc. Articles of Iron Carpentry Wood Misc. Manufactured Wood Articles	Misc. Rotating Electric Plant Parts Control Instruments of Gas or Liquid Valves	
F10: Materials II	F11: Automobile	F12: Chemicals I	F13: Chemicals II	F14: Machinery II	F15: Miscellaneous
Improved Wood Mineral Wool Central Heating Equipment	Vehicles Parts - Accessories Cars Iron Wire	Synthetic Rubber Acrylic Polymers Silicones	Aldehyde, Ketone Glycosides, Vaccines Medicaments	Parts of Metalworking Machine Tools Interchangeable Tool Parts Polishing Stones	Misc. Pumps Ash and Residues Chemical Wood Pulp of sulphite

# Results

## Interpretability

<b>Top Products (decay 30%)</b>	$B_{kd}$
Bovine	0.49
Miscellaneous Refrigeration Equipment	0.43
Radioactive Chemicals	0.41
Blocks of Iron and Steel	0.41
Rape Seeds	0.40
Animal meat, misc	0.39
Refined Sugars	0.38
Miscellaneous Tire Parts	0.38
Leather Accessories	0.38
Liquor	0.38
Bovine meat	0.38
Embroidery	0.37
Unmilled Barley	0.37
Dried Vegetables	0.36
Textile Fabrics Clothing Accessories	0.36
Horse Meat	0.35
Iron Bars and Rods	0.35
Analog Navigation Devices	0.35

(a) SVD

<b>Top Products (decay 30%)</b>	$B_{kd}$
Miscellaneous Animal Oils	0.78
Bovine and Equine Entrails	0.72
Bovine meat	0.68
Preserved Milk	0.63
Equine	0.62
Butter	0.58
Misc. Animal Origin Materials	0.57
Glues	0.56

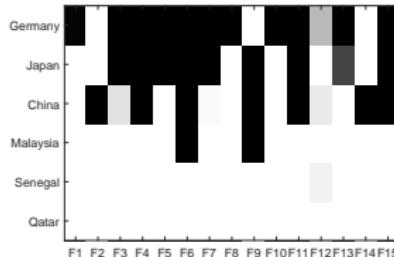
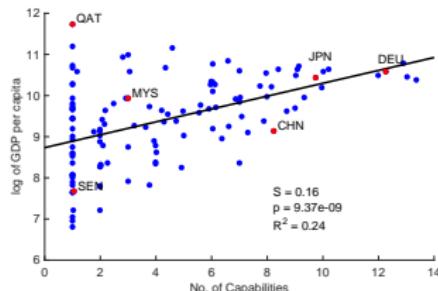
(b) S3R-IBP

# Results

## Interpretability

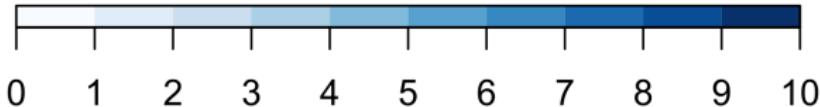
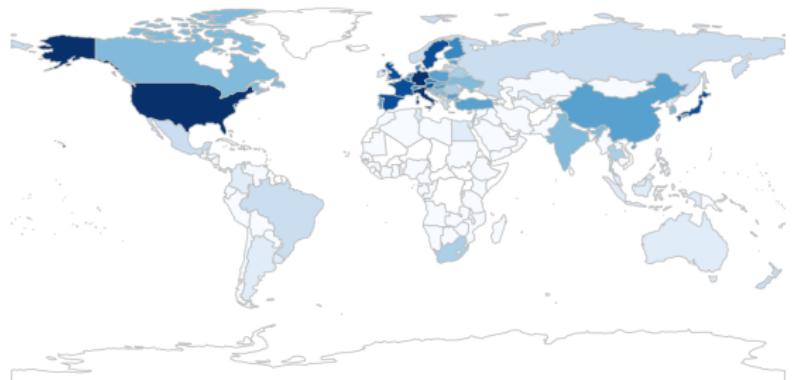
### Countries in Capability Space

- France = Belgium + Industrial Machinery
- Germany - Chemical = Austria
- Malaysia (Electronics) + Clothing → Phillipines
- Phillipines + Basic Processing → Indonesia, Vietnam



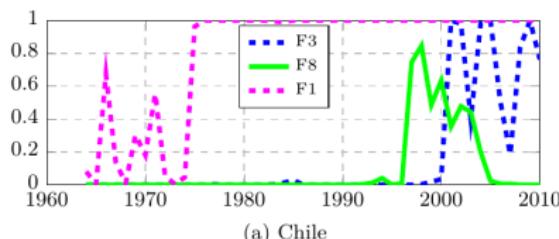
# Deep S3R-IBP: using a 2nd layer

- ① “Simple” and “advanced” capabilities
- ② Countries divided in two big groups: “quiescence” trap.

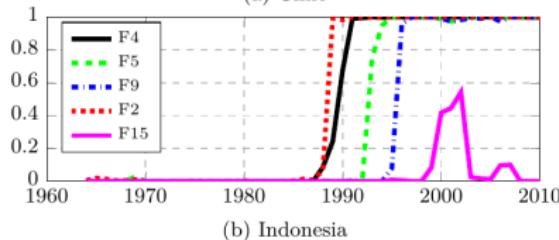


# Temporal Dynamics

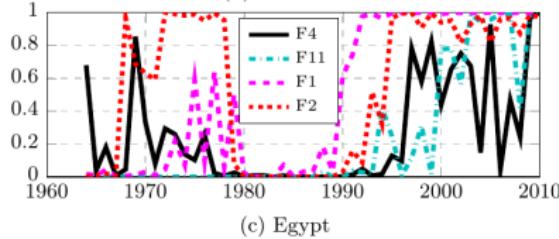
Capabilities	
F0	Bias
F1	Agriculture
F2	Clothing I
F3	Farming
F4	Clothing II
F5	Electronics I
F6	Processed Materials
F7	Electronics II
F8	Materials I
F9	Machinery I
F10	Materials II
F11	Automobile
F12	Chemicals I
F13	Chemicals II
F14	Machinery II
F15	Miscellaneous



(a) Chile



(b) Indonesia



(c) Egypt

# Conclusion

- ❶ BNP model for data exploration in high-dim count data.
- ❷ **interpretable** and **structured** solutions.
- ❸ Analysis of productive structure of world economies.

## Future works

- **Time-dependent extension** with Markovian activation of features and smooth variation of capabilities.

Thank you for listening! Any question?

[melanie@tsc.uc3m.es](mailto:melanie@tsc.uc3m.es)



# Sources and References

-  C. Bishop: *Pattern Recognition and Machine Learning*, 2006.
-  K. P. Murphy: *Machine Learning: a Probabilistic Perspective*, 2012.
-  D. J.C. MacKay: *Information Theory, Inference, and Learning Algorithms*, 2003.
-  Z. Ghahramani & C. E. Rasmussen, Slides for *Machine Learning Course* at Cambridge University.
-  S. J. Gershman, D.M. Blei: A tutorial on Bayesian nonparametric models, 2012.
-  Y.W. Teh: Slides for *Probabilistic and Bayesian Machine Learning*, UC3M, 2010.
-  M. N. Schmidt & M. Morup: *Advanced Topics in Machine Learning*, MLSS, DTU, 2013.
-  D. B. Dunson: *Nonparametric Bayes Applications to Biostatistics*, 2010.