

# Applications of latent variable models for data exploration and uncertainty quantification

March 15, 2019

Melanie F. Pradier



Center for Research on  
Computation and Society

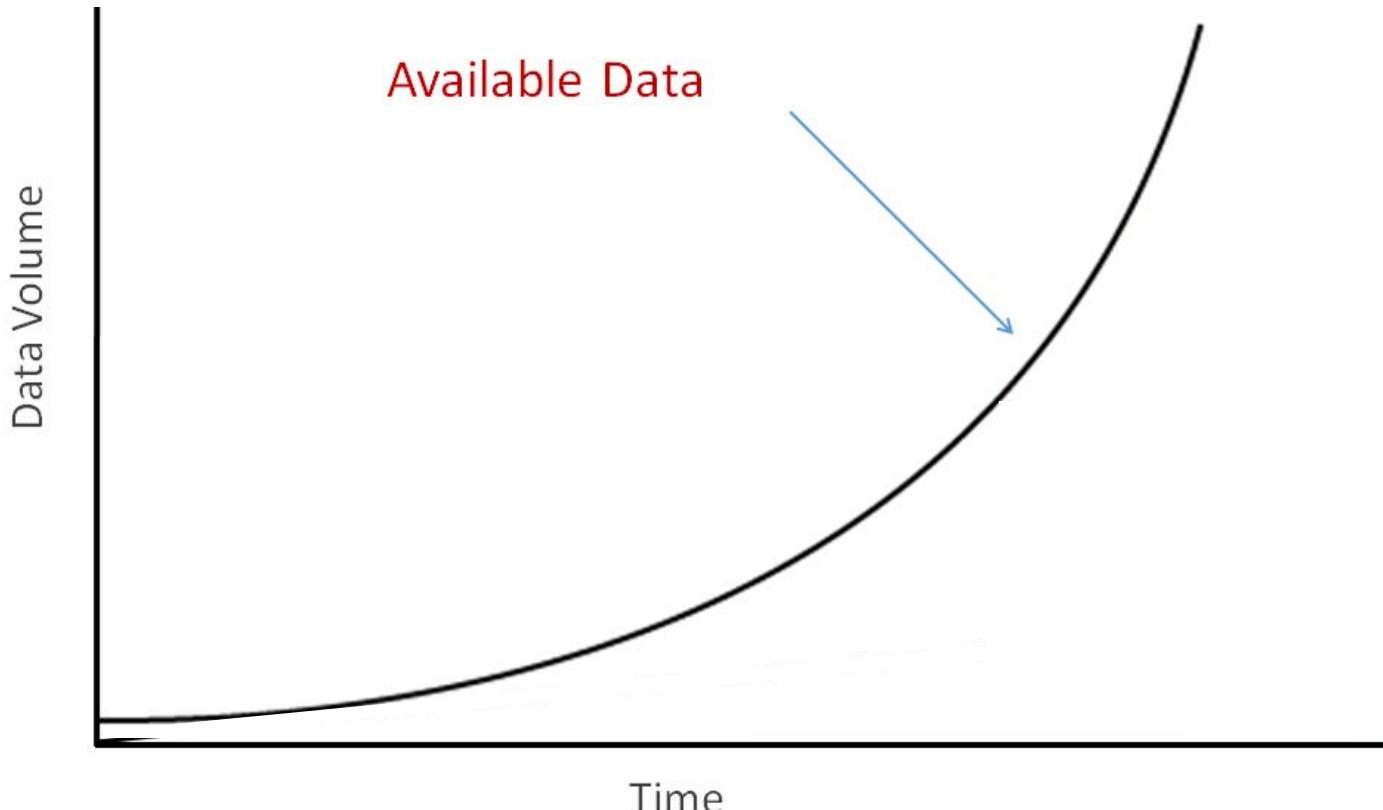
at Harvard John A. Paulson School of Engineering and Applied Sciences



**HDSI**

Harvard Data  
Science Initiative

# Data everywhere!



# Huge amount of opportunities...



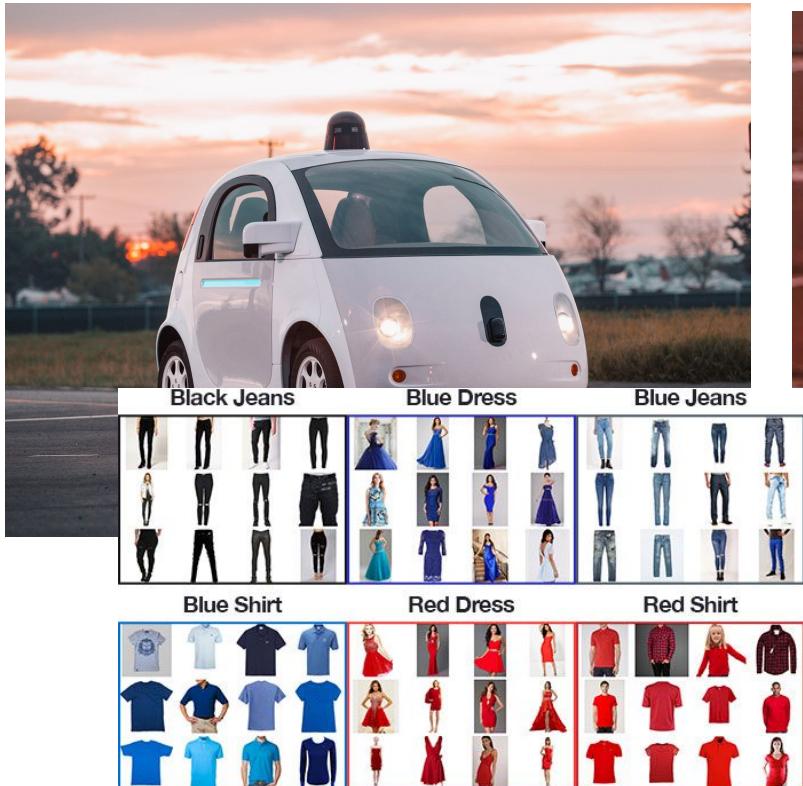
# Huge amount of opportunities...



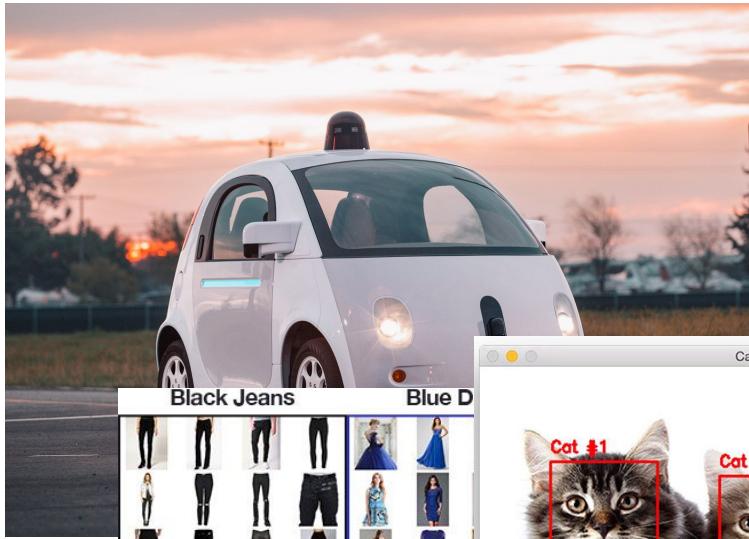
# Huge amount of opportunities...



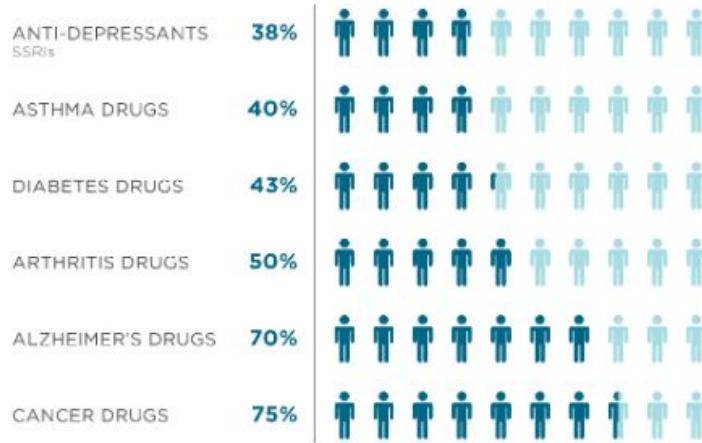
# Huge amount of opportunities...



# Huge amount of opportunities...



# ...but still many challenges



Prognostic

Is it likely to  
develop  
this cancer?

Diagnostic

What type of  
cancer is it?

Predictive

Is this the  
optimal  
drug for my  
cancer?

Pharmacodynamics

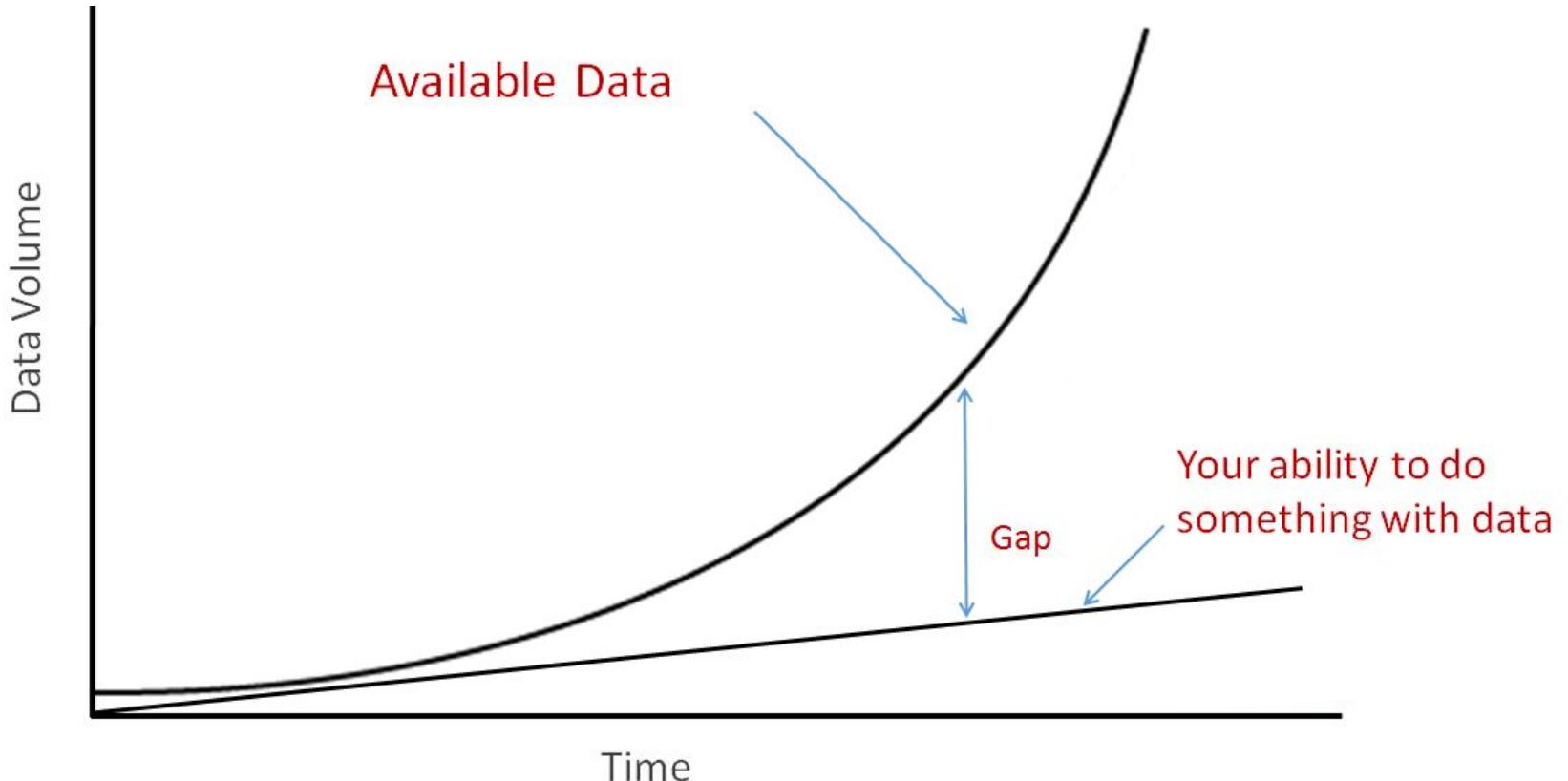
What's the  
optimal dose  
for my body?

Recurrence

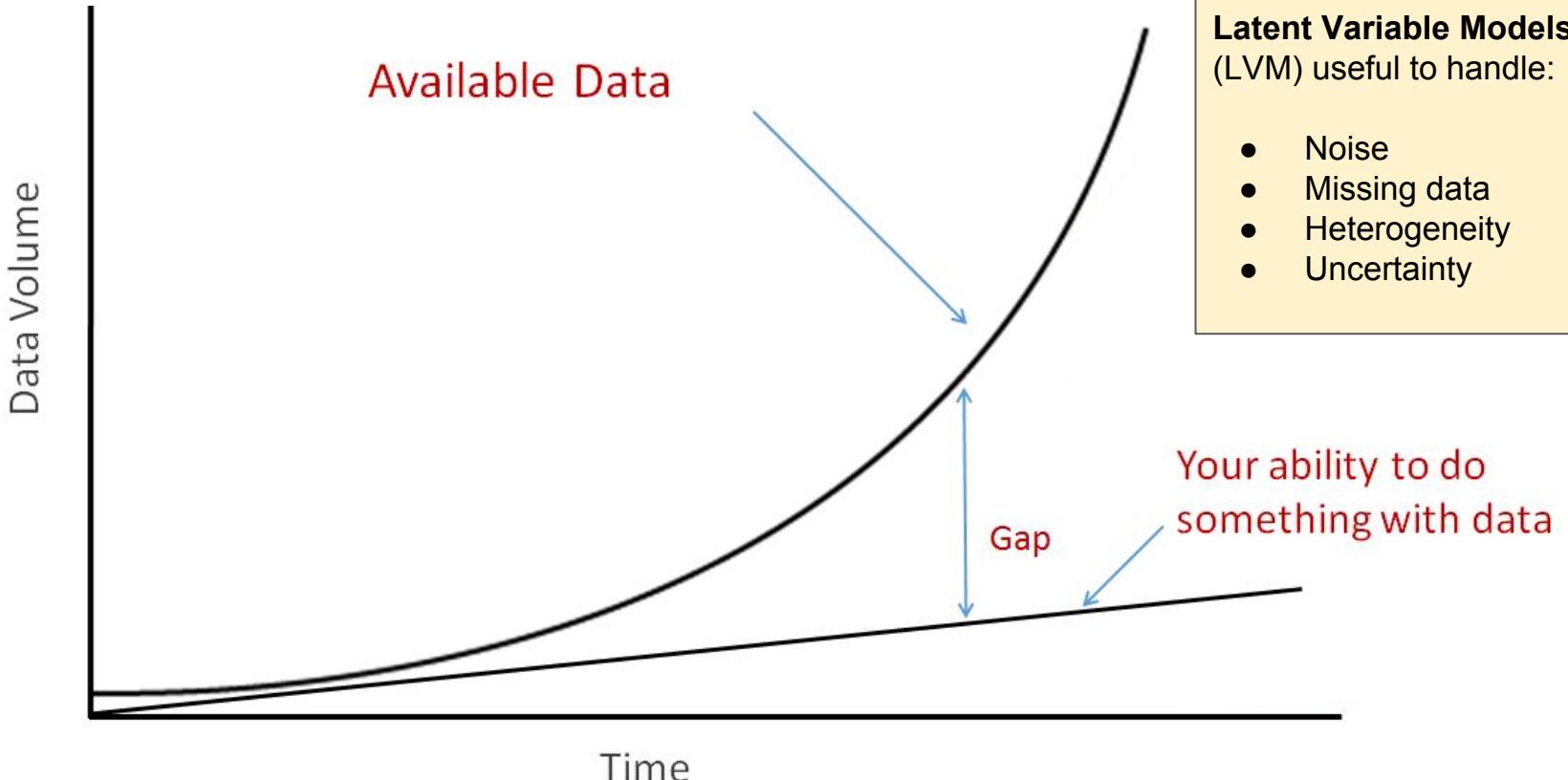
Will the  
cancer  
return?

Specially in high-stake decision scenarios

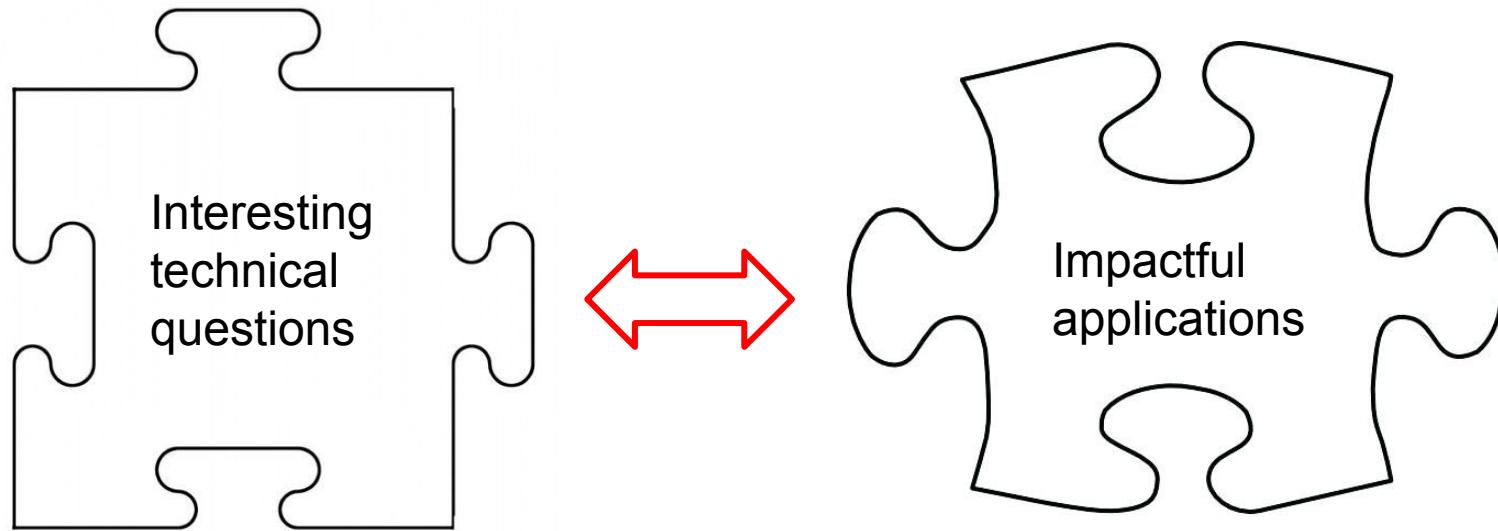
# Bridging the gap



# Bridging the gap



# My research: probabilistic models for societal needs



- Design probabilistic models (modeling/inference) for real-world applications
- Crucial: multidisciplinary collaboration

# Agenda from now on...

Applications of Latent Variable Models (LVMs) for:

1. Data Exploration
  - Biomarker discovery in clinical trials
2. Uncertainty Quantification
  - Inference framework for Bayesian neural networks

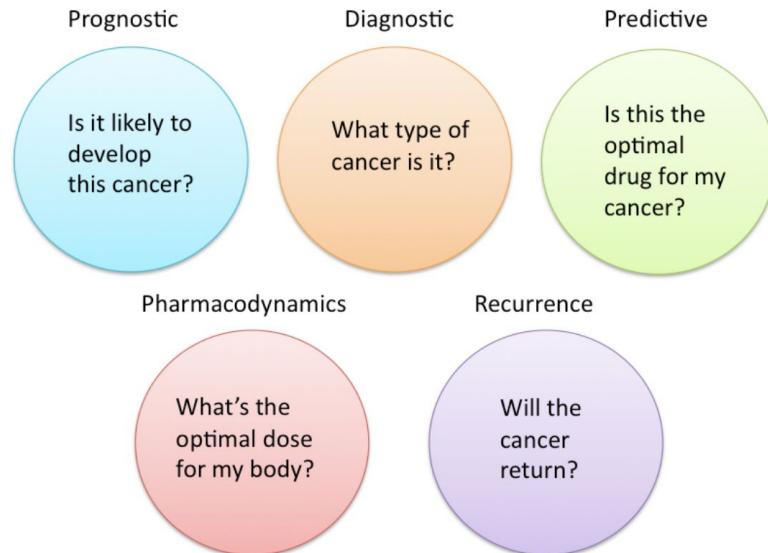
# Goal 1: Data exploration

## Objective: Biomarker discovery

Biomarkers used everywhere, e.g.,

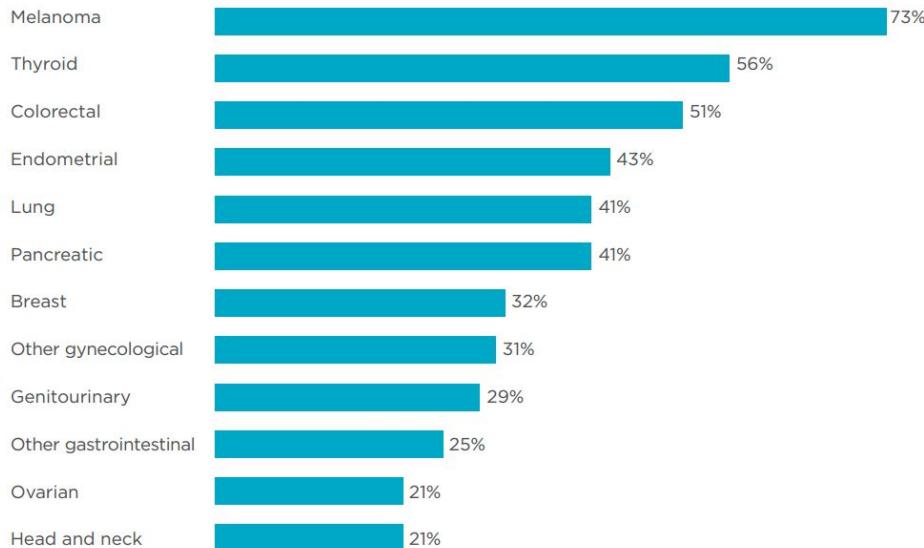
- Prostate-specific antigen (PSA) to diagnose prostate cancer
- Estrogen / progesterone to predict sensitivity to endocrine therapy in breast cancer
- KRAS mutation to predict resistance to EGFr antibody treatment

Biomarker = "any variable that can be used as an indicator of a particular disease state"



# Biomarker discovery is expensive

TACKLING TUMORS: Percentage of patients whose tumors were driven by certain genetic mutations that could be targets for specific drugs, by types of cancer.



Source: *Wall Street Journal* Copyright 2011 by DOW JONES & COMPANY, INC. Reproduced with permission of DOW JONES & COMPANY, INC.

## ANNUAL COST OF CANCER DRUGS

New cancer medicines now routinely cost more than \$100,000 yearly, which can create hardships even for insured patients. Top 10 oncological drugs by annual cost:

<b>Omacetaxine</b> for chronic myeloid leukemia	\$168,366
<b>Ibrutinib</b> mantle cell lymphoma	\$157,440
<b>Crizotinib</b> non-small-cell lung cancer	\$156,544
<b>Pomalidomide</b> multiple myeloma	\$150,408
<b>Regorafenib</b> colorectal cancer	\$141,372
<b>Sorafenib</b> papillary thyroid cancer	\$140,984
<b>Ponatinib</b> chronic myeloid leukemia <sup>1</sup>	\$137,952
<b>Trametinib</b> malignant melanoma	\$125,280
<b>Lenalidomide</b> mantle cell lymphoma	\$124,870
<b>Cabozantinib</b> medullary thyroid cancer	\$118,800

Among drugs approved between 2009 and 2013 by the Food and Drug Administration

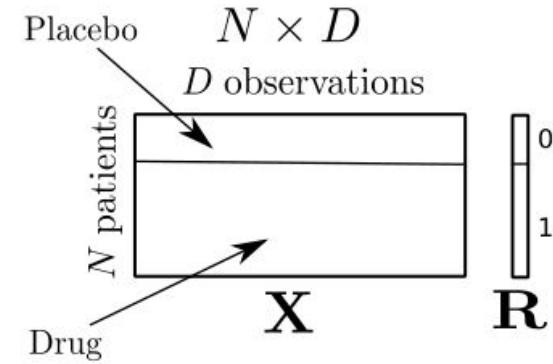
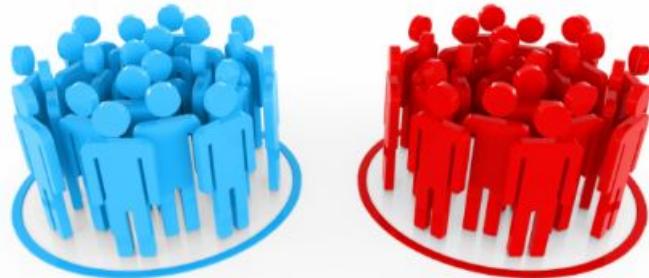
<sup>1</sup>—Also for Ph+ acute lymphoblastic leukemia

SOURCE: JAMA Oncology, 2015

George Petras, USA TODAY

# Problem formulation

Clinical trial scenario

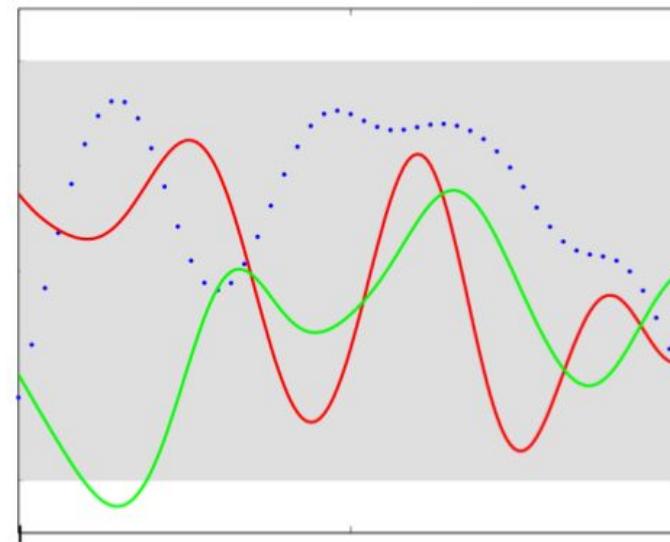
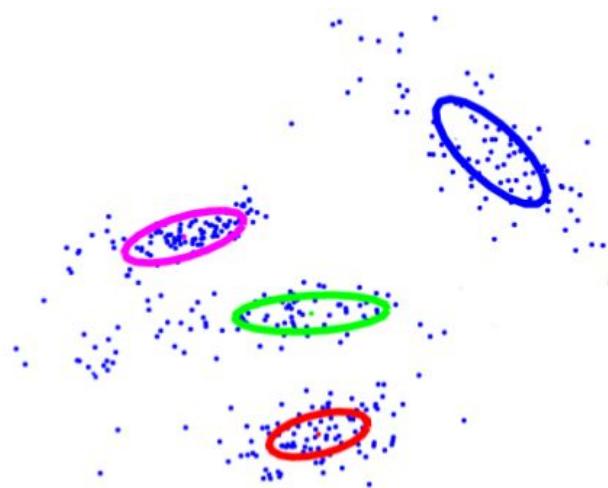


We want to discover:

- ① Indicators of disease progression: prognostic biomarkers
- ② Indicators of (positive) drug response: predictive biomarkers

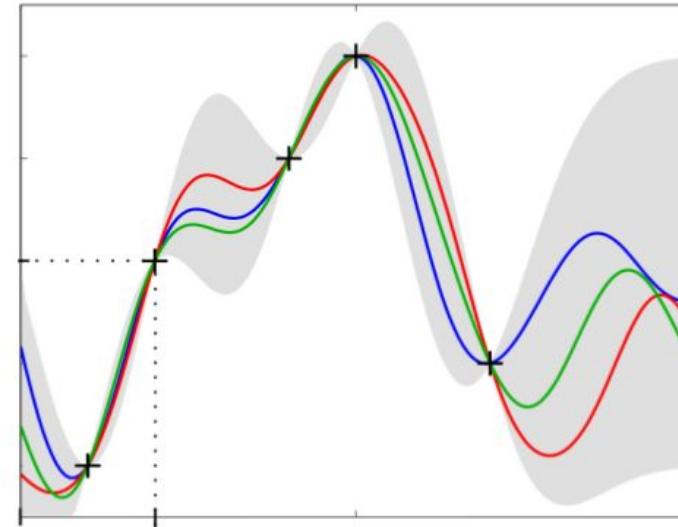
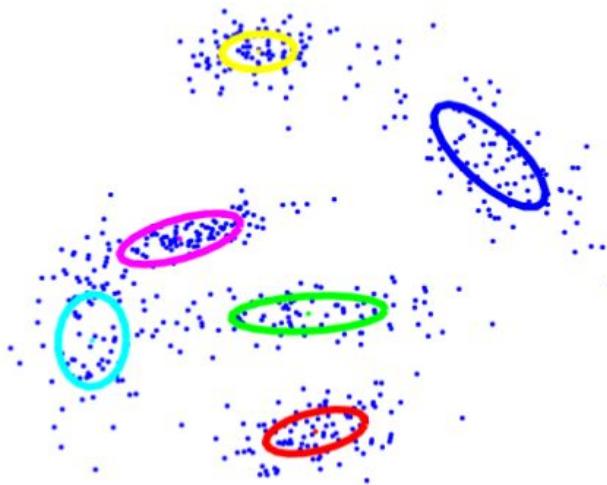
# Bayesian nonparametrics

- Bayesian: to handle uncertainty  $p(\Phi|\mathbf{X}) \propto p(\mathbf{X}|\Phi)p(\Phi)$
- Nonparametric: to adapt model complexity depending on input data (hypothesis generation)

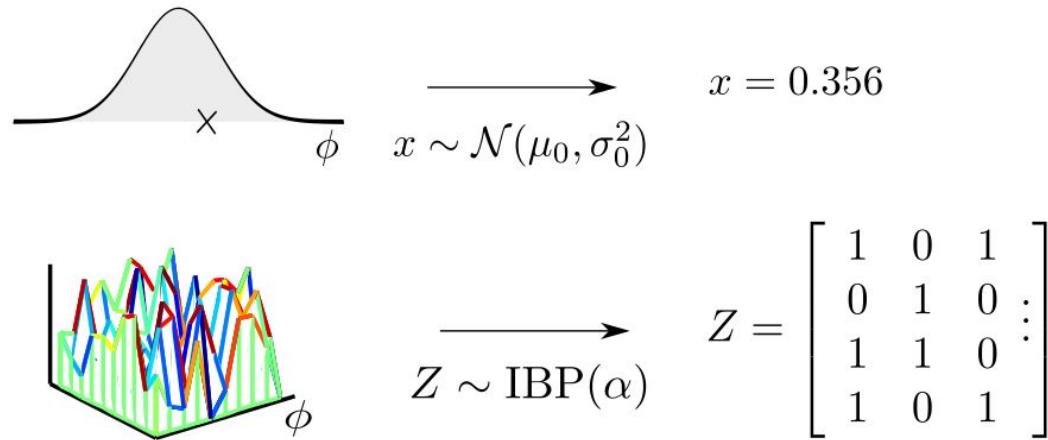


# Bayesian nonparametrics

- Bayesian: to handle uncertainty  $p(\Phi|X) \propto p(X|\Phi)p(\Phi)$
- Nonparametric: to adapt model complexity depending on input data (hypothesis generation)



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



- Prior over binary matrices with infinite number of columns
- Rows  $\equiv$  observations; columns  $\equiv$  features
- $Z \sim \text{IBP}(\alpha)$
- $\alpha$ : concentration parameter

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

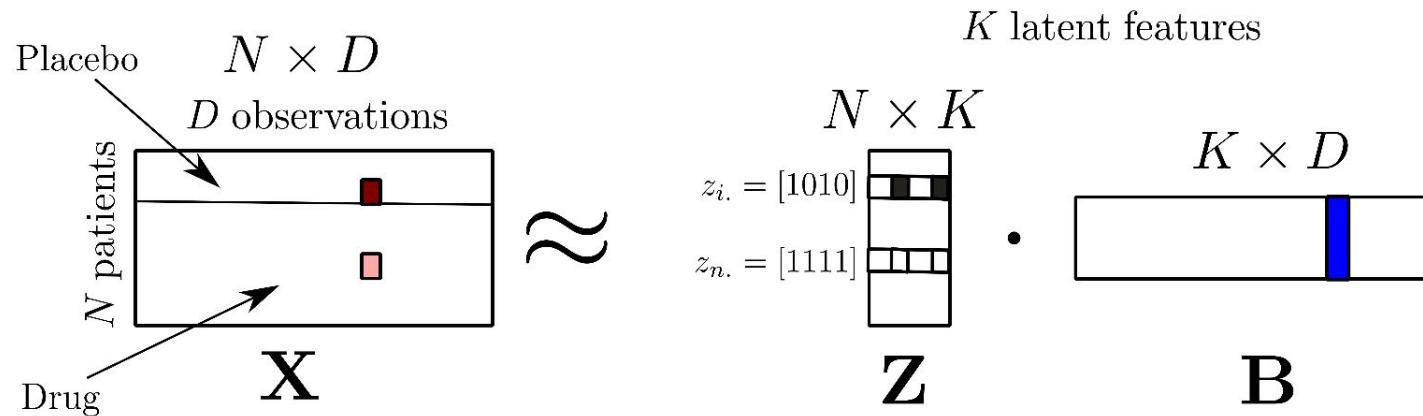


# Indian Buffet Process (Ghahramani et.al, 2006)

	1	2	3	4	5	6	...
1	1	1	1	0	0	0	
2	1	0	1	1	0	0	
3	0	1	1	0	1	1	
:	:						

The diagram illustrates the Indian Buffet Process (IBP) through a matrix representation. At the top, six different types of Indian cuisine are shown in small images. Below the images is a horizontal ellipsis (three dots). To the left of the matrix, there are three user icons: a man with brown hair in a blue shirt, a woman with long dark hair in a grey shirt, and a man with blonde hair in a grey suit. Below these icons is a vertical ellipsis (two dots). The matrix itself has 3 rows (users) and 6 columns (food items). The values in the matrix represent binary variables indicating whether each user has tried or not tried each food item. The first row shows users 1, 2, and 3 having tried all three items, while users 4, 5, and 6 have not. The second row shows users 1 and 3 having tried item 2, while users 2, 4, 5, and 6 have not. The third row shows users 2 and 3 having tried item 1, while users 1, 4, 5, and 6 have not.

# Infinite latent feature model (intuition)

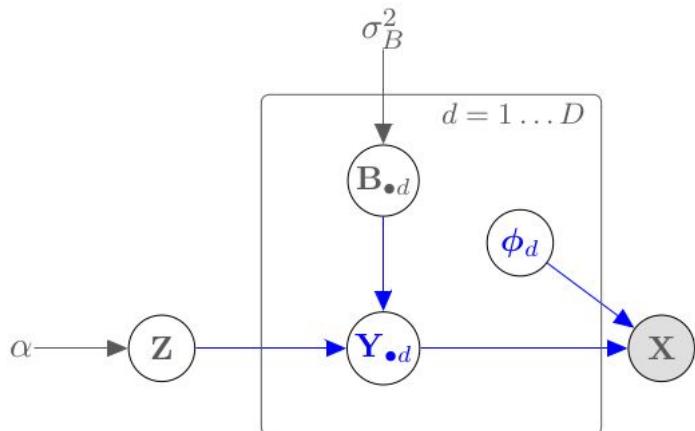


- $x_{id} = 173 \text{ ml/dL} = 73 + 0 + 100 \text{ ml/dL}$
- $x_{nd} = 136 \text{ ml/dL} = 86 + 40 + 60 - 50 \text{ ml/dL}$

# General Latent Feature Model (GLFM)

Latent feature model for heterogeneous datasets

- Link functions  $T_d$  depend on type of data for each dimension  $d$

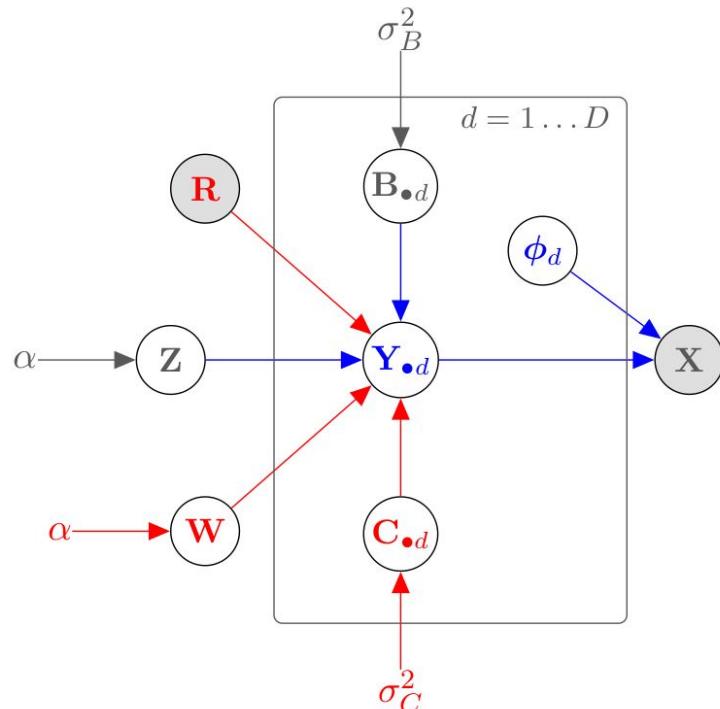


$$\begin{aligned}x_{nd} &= T_d(y_{nd}; \phi_d) \\y_{nd} | \mathbf{Z}, \mathbf{B} &\sim \mathcal{N}(\mathbf{Z}_{n\bullet} \mathbf{B}_{\bullet d}, \sigma_y^2) \\B_{kd} &\sim \mathcal{N}(0, \sigma_B^2) \\\mathbf{Z} &\sim \text{IBP}(\alpha)\end{aligned}$$

Open-source python code

<https://github.com/ivaleraM/GLFM>

# Case-Control Indian Buffet Process (C-IBP)



$R_n$ : drug indicator por patient  $n$

$$x_{nd} = T_d(y_{nd}; \phi_d)$$

$$y_{nd}|Z, W, B, C, R \sim$$

$$\mathcal{N}(Z_n \bullet B_{\bullet d} + \mathbb{1}[R_n = 1] W_n \bullet C_{\bullet d}, \sigma_y^2)$$

$$B_{kd} \sim \mathcal{N}(0, \sigma_B^2)$$

$$Z \sim \text{IBP}(\alpha)$$

$$C_{kd} \sim \mathcal{N}(0, \sigma_C^2)$$

$$W \sim \text{IBP}(\alpha)$$

- **Inference:** MCMC approach with accelerated Gibbs sampling
- **Biomarker discovery:** statistical multiple hypothesis testing

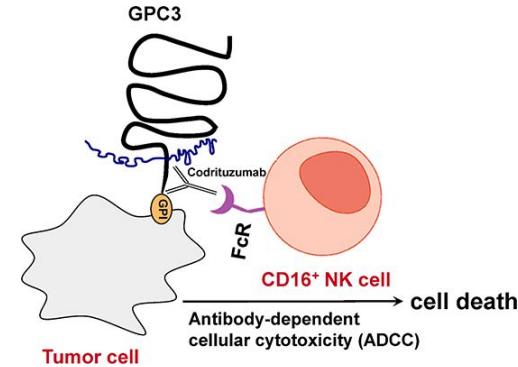
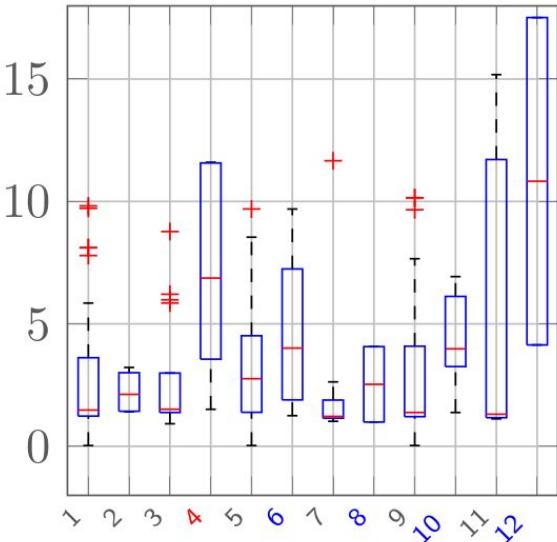
# Results: subpopulations

GPC3 Antibody Treatment against Liver Cancer (J. Hepatology. 2016 Apr, Abou-Alfa et.al.)

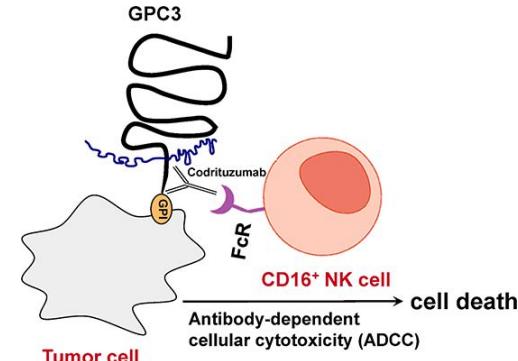
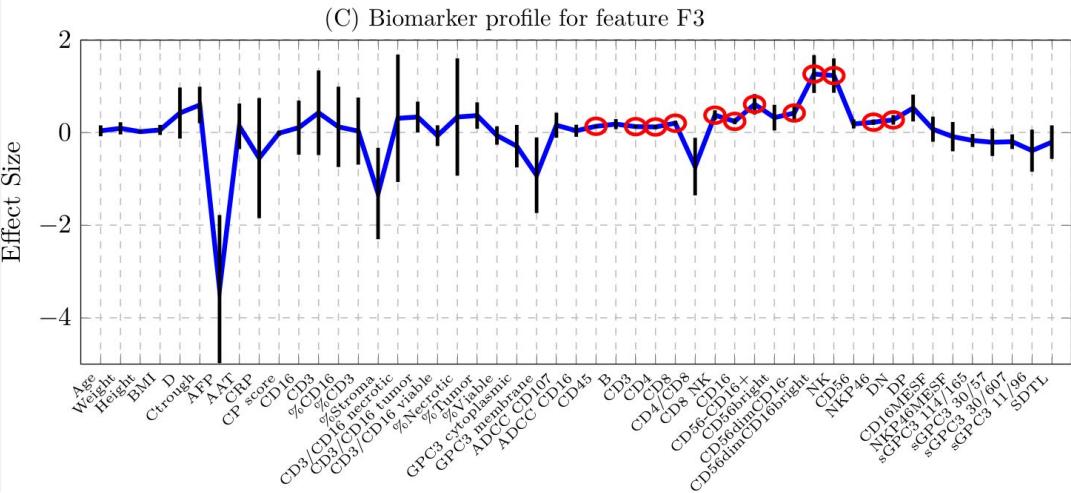
- 180 patients: 60 took a placebo, 120 took the drug
- PFS: Progression Free Survival

Sub-population	Drug Identifier				Size (number of patients)	Mean PFS (months)	Median PFS (months)
		F1	F2	F3			
1.	0	0	0	0	33.37	3.06	1.65
2.	0	0	1	0	4.07	2.29	2.24
3.	0	1	0	0	17.84	2.72	1.81
4.	0	1	1	0	4.72	7.05	7.18
5.	1	0	0	0	51.52	3.22	2.55
6.	1	0	0	1	16.77	4.17	3.65
7.	1	0	1	0	8.38	1.74	1.33
8.	1	0	1	1	2.07	2.69	2.65
9.	1	1	0	0	29.88	3.36	2.03
10.	1	1	0	1	4.90	4.44	4.34
11.	1	1	1	0	4.53	6.31	5.31
12.	1	1	1	1	1.94	10.04	10.01

PFS (months)



## Results: biomarker profiles

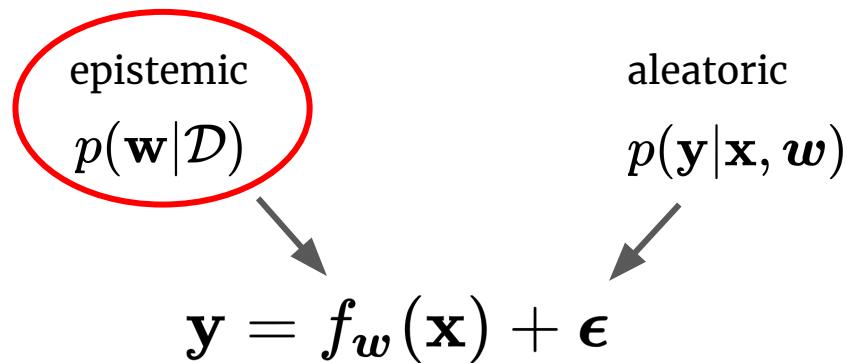


## **Take-away message:**

- LVMs useful to identify hidden patterns underlying data
  - Challenge addressed: data heterogeneity (both across dimensions and observations)

# Goal 2: Uncertainty quantification

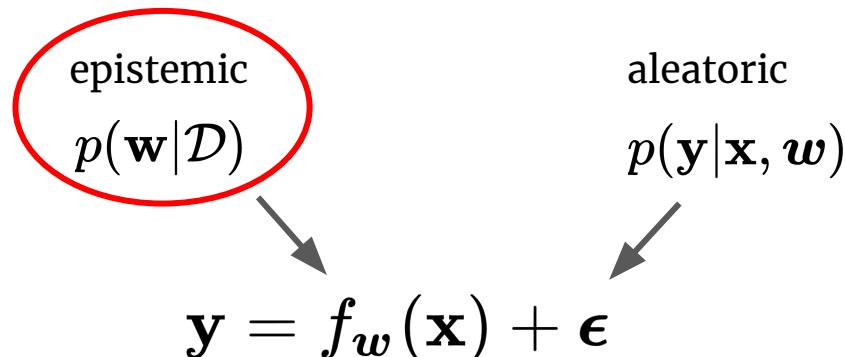
Two sources of uncertainty



[Depeweg et.al, 2017]

# Goal 2: Uncertainty quantification

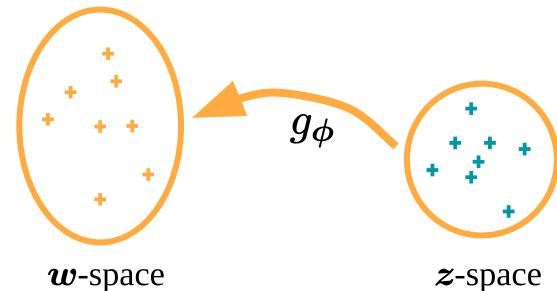
Two sources of uncertainty



[Depeweg et.al, 2017]

High-level idea

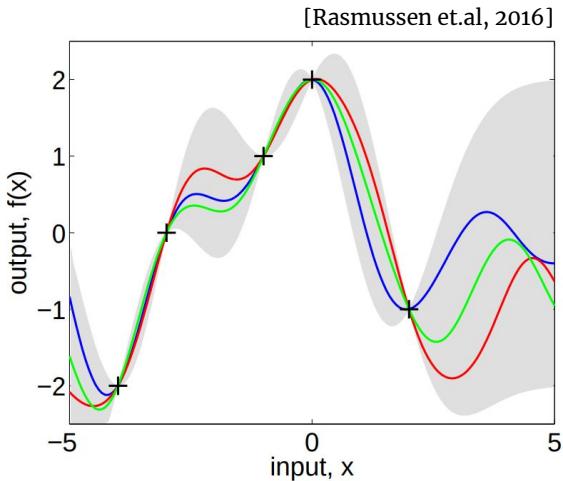
- Approximate  $f_{\mathbf{w}}$  with a Bayesian Neural Network



- Modeling + inference contributions

# How to estimate function uncertainty?

## Gaussian Process (GP)



$$f(x) \sim \text{GP} (m(x), k(x, x'))$$

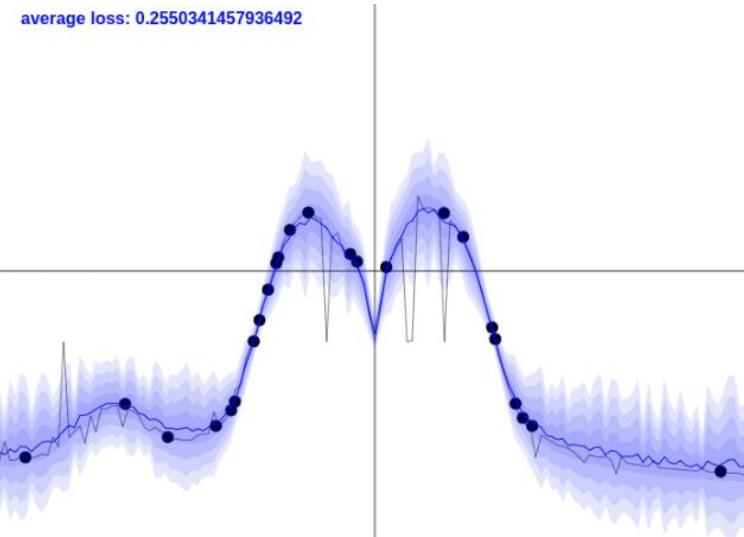
## Drawbacks of GPs

- Scalability
- Kernel learning is not trivial

## Alternative: Neural Networks with uncertainty

- Ensemble of Neural Networks  
[Lakshminarayanan et al., 2017; Pearce et.al, 2018]
- Bayesian Neural Networks  
[Buntine et al., 1991; MacKay, 1992; Neal, 1993]

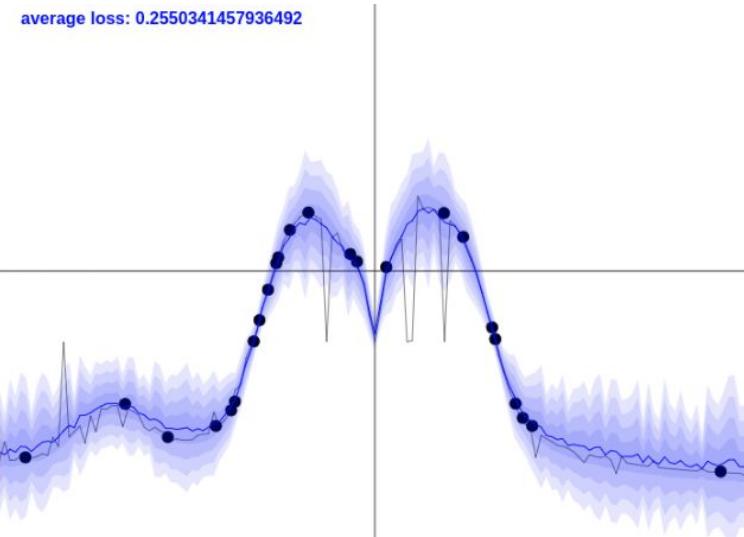
# Bayesian Neural Network (BNN)



$$\mathbf{y} = f_{\mathbf{w}}(\mathbf{x}) + \boldsymbol{\epsilon} \quad \mathcal{D} = \{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^N$$
$$\mathbf{w} \sim \mathcal{N}(0, \sigma_w^2 \mathbf{I}), \quad \boldsymbol{\epsilon} \sim \mathcal{N}(0, \sigma_\epsilon^2 \mathbf{I})$$

[What my deep model does not know, post of Yarin Gal, 2015]

# Bayesian Neural Network (BNN)



[What my deep model does not know, post of Yarin Gal, 2015]

$$\mathbf{y} = f_{\mathbf{w}}(\mathbf{x}) + \boldsymbol{\epsilon} \quad \mathcal{D} = \{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^N$$

$$\mathbf{w} \sim \mathcal{N}(0, \sigma_w^2 \mathbf{I}), \quad \boldsymbol{\epsilon} \sim \mathcal{N}(0, \sigma_\epsilon^2 \mathbf{I})$$

Quantities of interest:

- Posterior of the weights  $p(\mathbf{w}|\mathcal{D})$
- Predictive distribution

$$p(\mathbf{y}^* | \mathbf{x}^*, \mathcal{D}) = \int p(\mathbf{y}^* | \mathbf{x}^*, \mathbf{w}) p(\mathbf{w}|\mathcal{D}) d\mathbf{w}$$

$$p(w|\mathcal{D})$$

is intractable!

Inference options:

- **Markov Chain Monte Carlo**  
Hamiltonian Monte Carlo [Neal, 1993]
- **Variational Inference**  
[Graves, 1993] [Blundell et.al, 2015]

# Variational Inference for BNNs

[Blundell et.al, 2015]

Objective: approximate  $p(\mathbf{w}|\mathcal{D})$

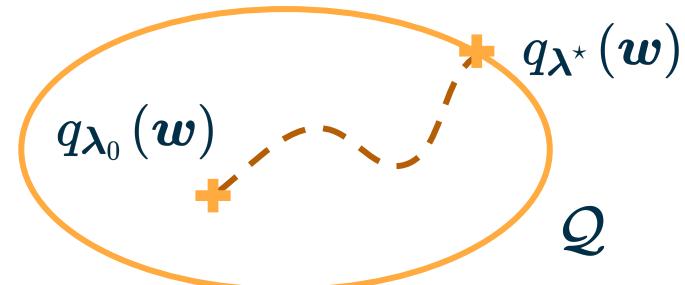
$$+ p(\mathbf{w}|\mathcal{D})$$

$$q_{\lambda}(\mathbf{w}) \in \mathcal{Q}$$

$$\underset{\lambda^*}{\operatorname{argmin}} D_{\text{KL}}\left(q_{\lambda}(\mathbf{w})||p(\mathbf{w}|\mathcal{D})\right)$$

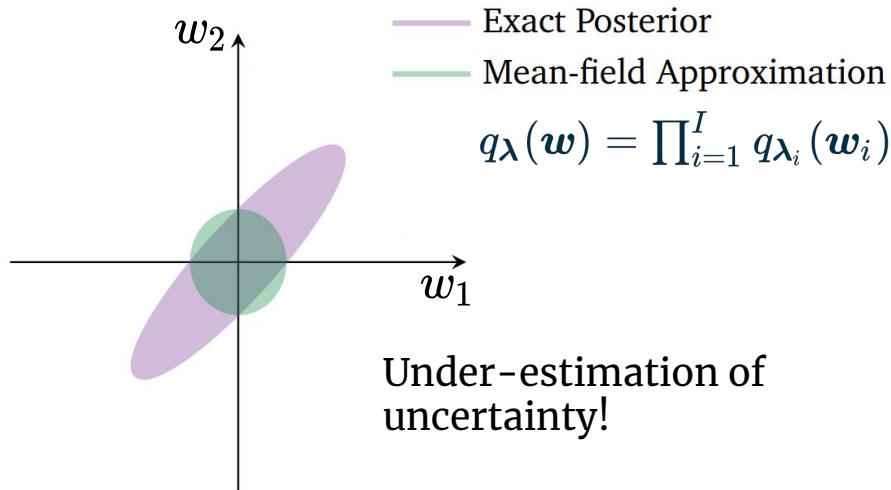


$$\underset{\lambda^*}{\operatorname{argmax}} \mathcal{L}(\lambda) = \mathbb{E}_q \left[ \log p(\mathbf{y}|\mathbf{x}, \mathbf{w}) \right] - D_{\text{KL}}(q_{\lambda}(\mathbf{w})||p(\mathbf{w}))$$



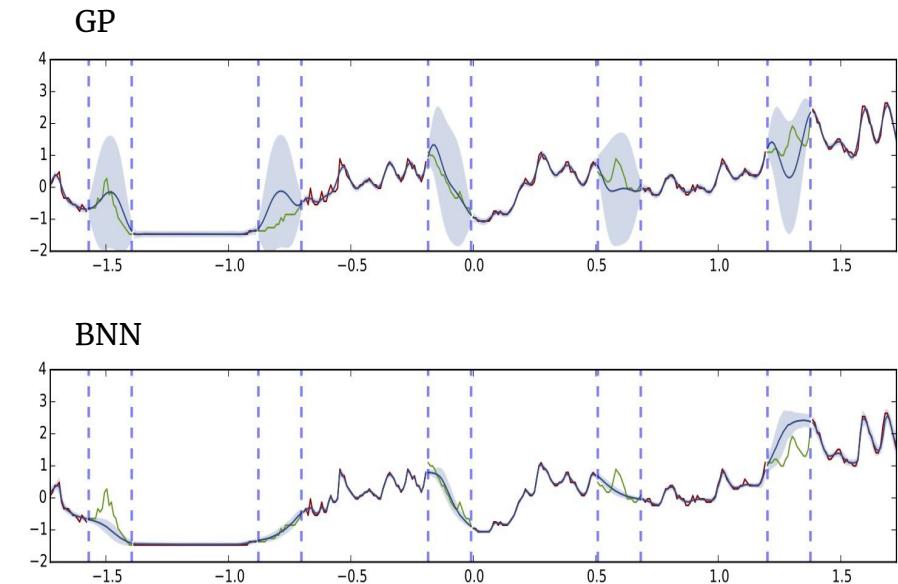
Black-box VI [Ranganath et.al, 2013] + reparametrization trick [Kingma et.al, 2014; Rezende et.al, 2015]

# Is mean-field VI good enough?



- More flexible variational approx. in weight Space [Louizos et al, 2016; 2017]

Example on solar irradiance dataset [Gal et.al, 2015]



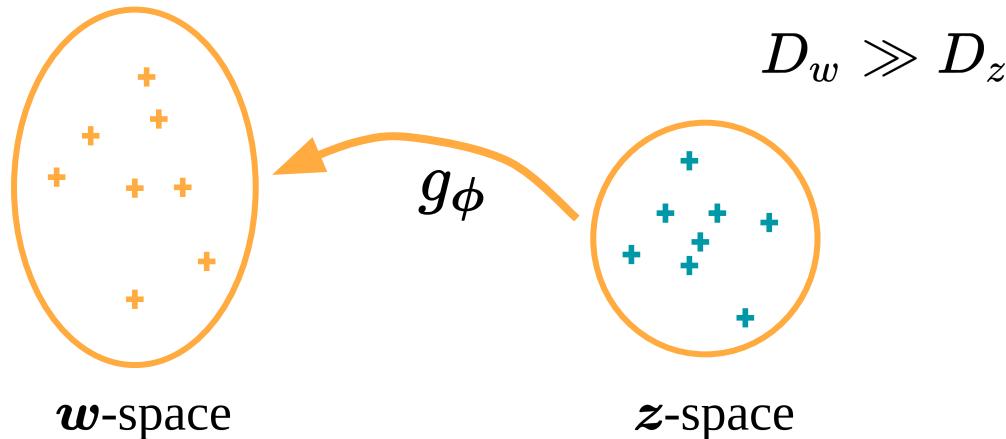
# Standard BNN

$$\begin{aligned}\mathbf{y} &= f_{\mathbf{w}}(\mathbf{x}) + \boldsymbol{\epsilon}, \quad \mathbf{w} \sim \mathcal{N}(0, \sigma_w^2 \mathbf{I}), \\ \boldsymbol{\epsilon} &\sim \mathcal{N}(0, \sigma_\epsilon^2 \mathbf{I})\end{aligned}$$

Weight redundancy  
[Denil et.al, 2013]

# Projected BNN

$$\begin{aligned} \mathbf{y} &= f_{\mathbf{w}}(\mathbf{x}) + \epsilon, \quad \mathbf{w} = g_{\phi}(\mathbf{z}), \quad \mathbf{z} \sim p(\mathbf{z}), \quad \phi \sim p(\phi), \\ \epsilon &\sim \mathcal{N}(0, \sigma_{\epsilon}^2 \mathbf{I}) \end{aligned}$$



# How about inference?

Objective: approximate  $p(\mathbf{w}|\mathcal{D})$

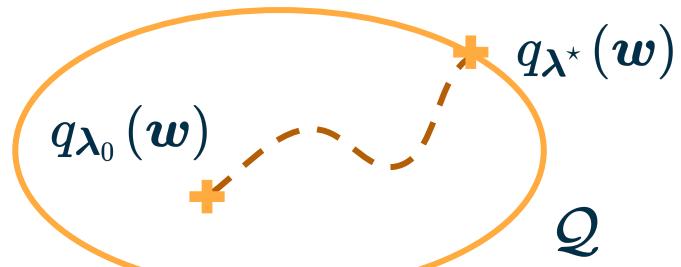
$$+ p(\mathbf{w}|\mathcal{D})$$

$$q_{\lambda}(\mathbf{w}) \in \mathcal{Q}$$

$$\underset{\lambda^*}{\operatorname{argmin}} D_{\text{KL}}(q_{\lambda}(\mathbf{w}) || p(\mathbf{w}|\mathcal{D}))$$



$$\underset{\lambda^*}{\operatorname{argmax}} \mathcal{L}(\lambda) = \mathbb{E}_q \left[ \log p(\mathbf{y}|\mathbf{x}, \mathbf{w}) \right] - D_{\text{KL}}(q_{\lambda}(\mathbf{w}) || p(\mathbf{w}))$$



Black-box VI [Ranganath et.al, 2013] + reparametrization trick [Kingma et.al, 2014; Rezende et.al, 2015]

# How about inference?

Objective: approximate  $p(\mathbf{z}, \boldsymbol{\phi} | \mathcal{D})$

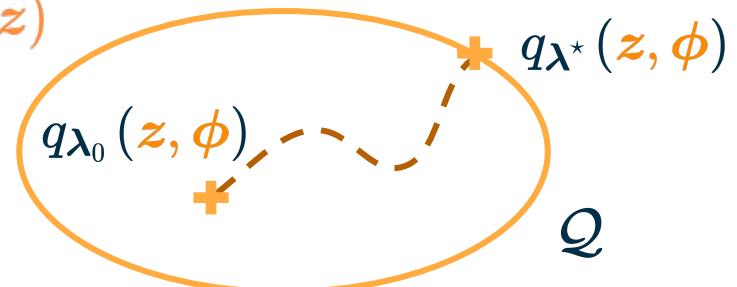
$$+ p(\mathbf{z}, \boldsymbol{\phi} | \mathcal{D})$$

$$\mathbf{z} \sim q_{\boldsymbol{\lambda}_z}(\mathbf{z}), \quad \boldsymbol{\phi} \sim q_{\boldsymbol{\lambda}_{\boldsymbol{\phi}}}(\boldsymbol{\phi}), \quad \mathbf{w} = g_{\boldsymbol{\phi}}(\mathbf{z})$$

$$\operatorname{argmin}_{\boldsymbol{\lambda}^*} D_{\text{KL}}(q_{\boldsymbol{\lambda}}(\mathbf{z}, \boldsymbol{\phi}) || p(\mathbf{z}, \boldsymbol{\phi} | \mathcal{D}))$$



$$\operatorname{argmax}_{\boldsymbol{\lambda}^*} \mathcal{L}(\boldsymbol{\lambda}) = \mathbb{E}_q \left[ \log p(\mathbf{y} | \mathbf{x}, g_{\boldsymbol{\phi}}(\mathbf{z})) \right] - D_{\text{KL}}(q_{\boldsymbol{\lambda}_z}(\mathbf{z}) || p(\mathbf{z})) - D_{\text{KL}}(q_{\boldsymbol{\lambda}_{\boldsymbol{\phi}}}(\boldsymbol{\phi}) || p(\boldsymbol{\phi}))$$



Black-box VI [Ranganath et.al, 2013] + reparametrization trick [Kingma et.al, 2014; Rezende et.al, 2015]

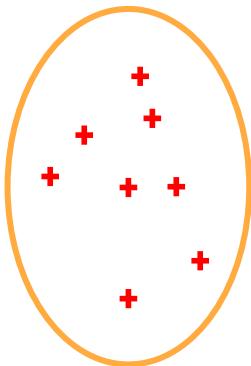
$$\underset{\lambda^*}{\operatorname{argmin}} D_{\text{KL}} \left( q_{\lambda}(z, \phi) || p(z, \phi | \mathcal{D}) \right)$$

jointly is hard!

Our solution: find smart initialization

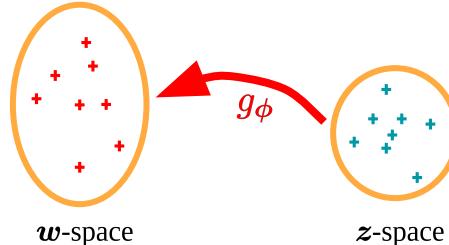
# Solution: 3-stage Inference Framework

1. Characterize weight space



Train ensemble of  
neural networks

2. Find point estimate  $g_\phi$



Train an  
autoencoder

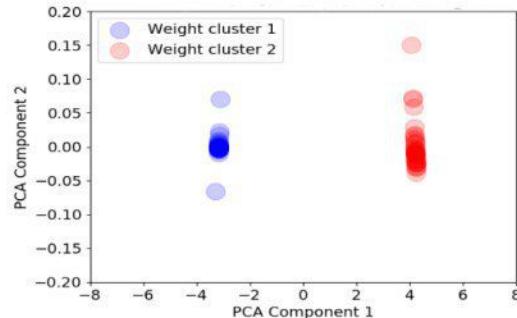
3. Black-box VI (BBVI)

$$D_{\text{KL}} \left( q_\lambda(z, \phi) || p(z, \phi | \mathcal{D}) \right)$$

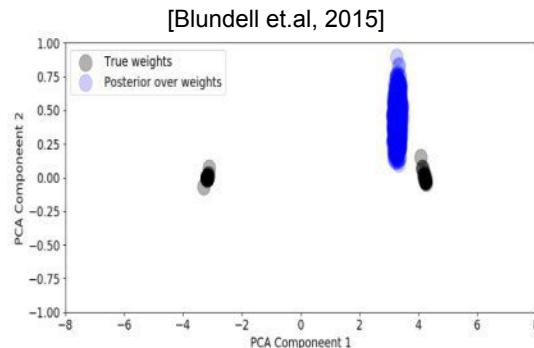
Principled BBVI with  
smart initialization

# Results

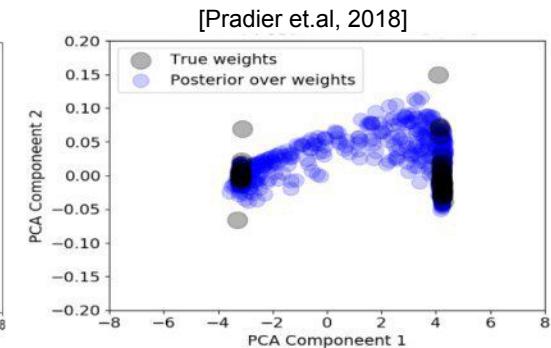
# Illustrative Toy Example



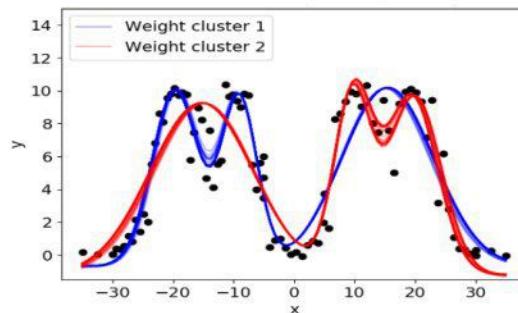
(a) Projection of true weights



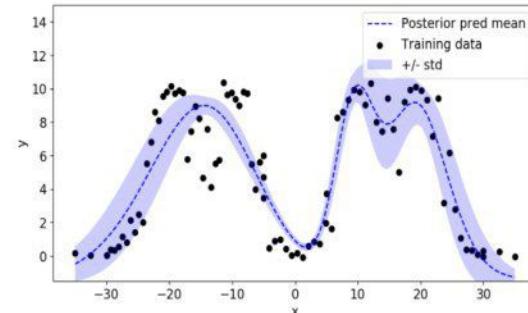
(b) BbB posterior over weights



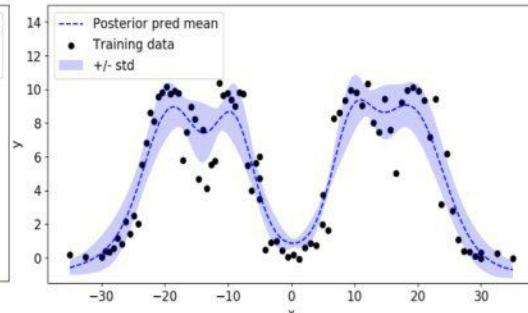
(c) Proj-BNN posterior over weights



(d) Functions from true weights

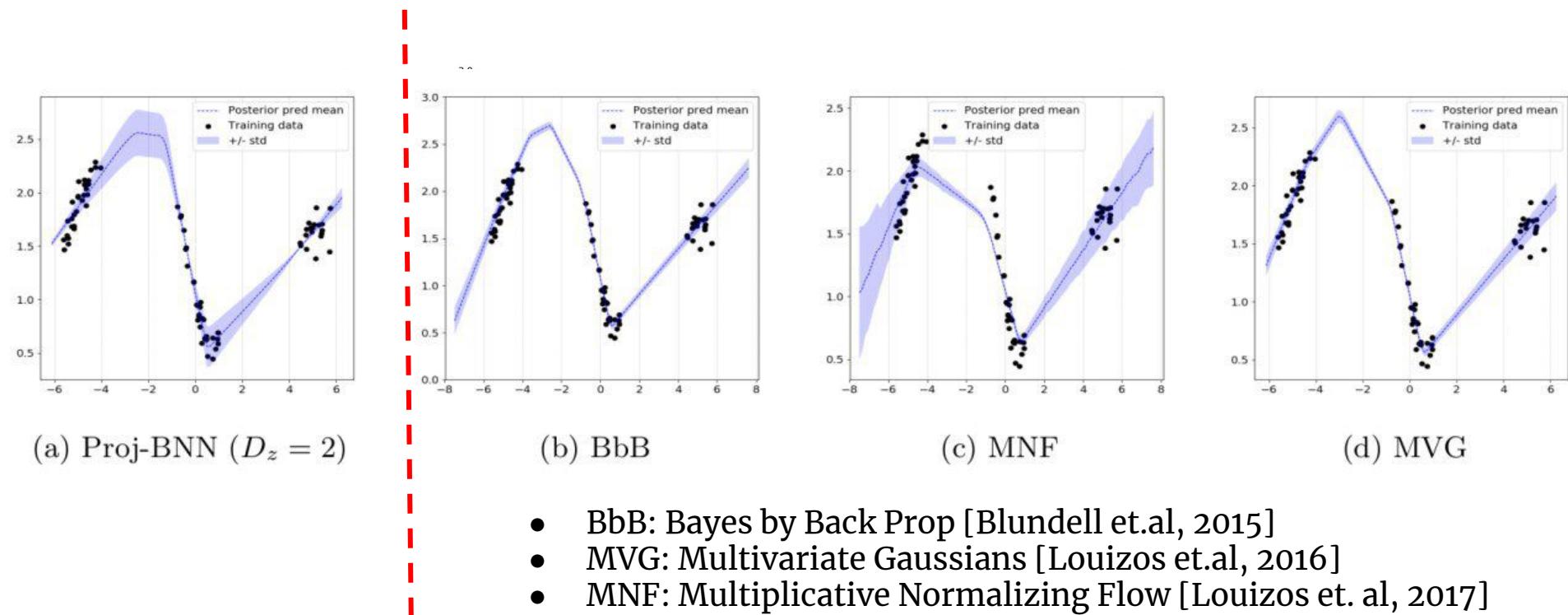


(e) BbB posterior predictive

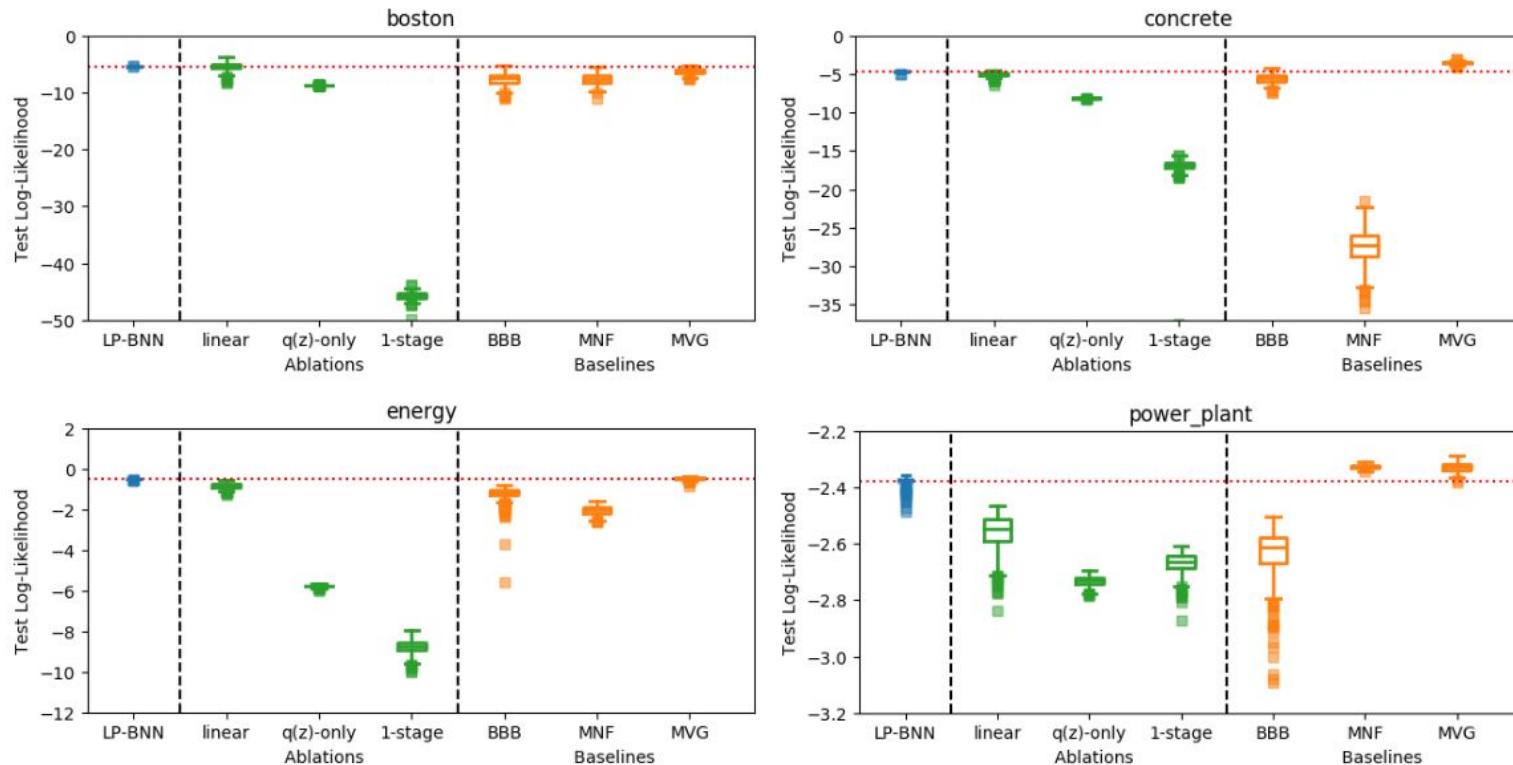


(f) Proj-BNN posterior predictive

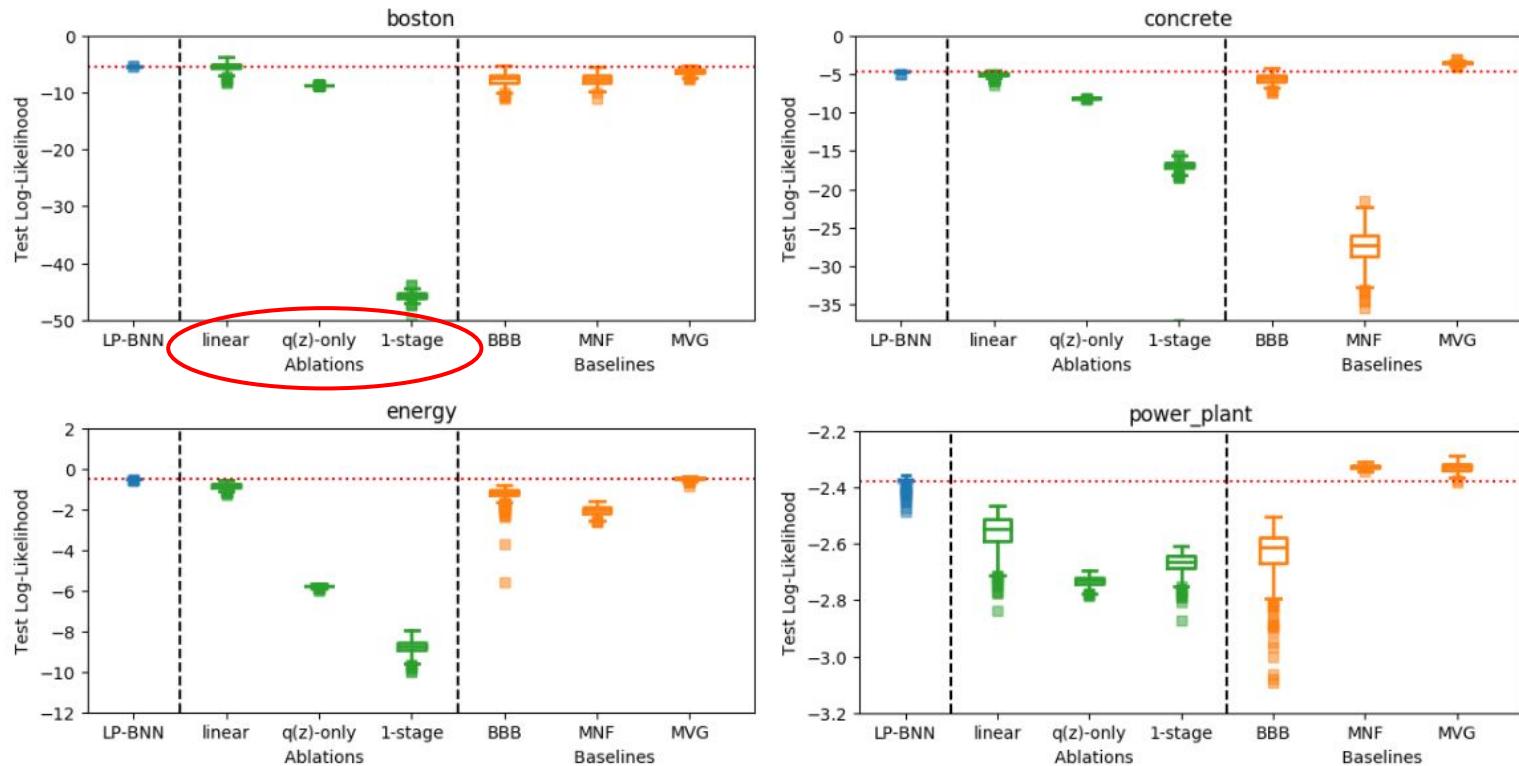
# Results: Uncertainty estimation



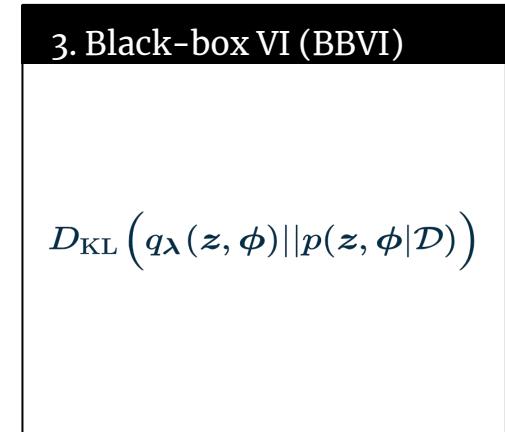
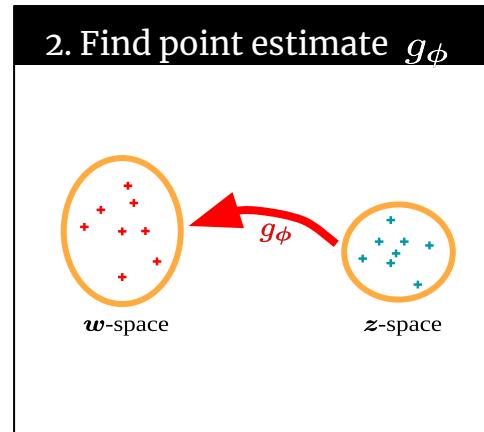
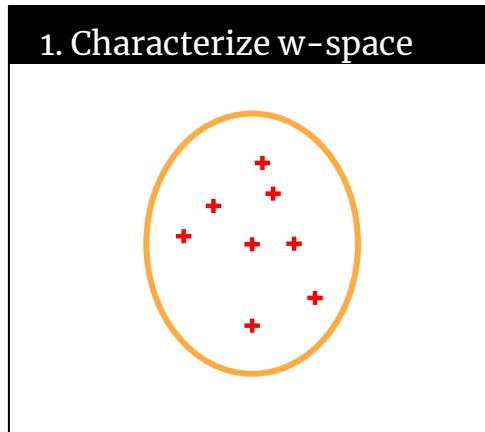
# Results: Generalization



# Results: Generalization



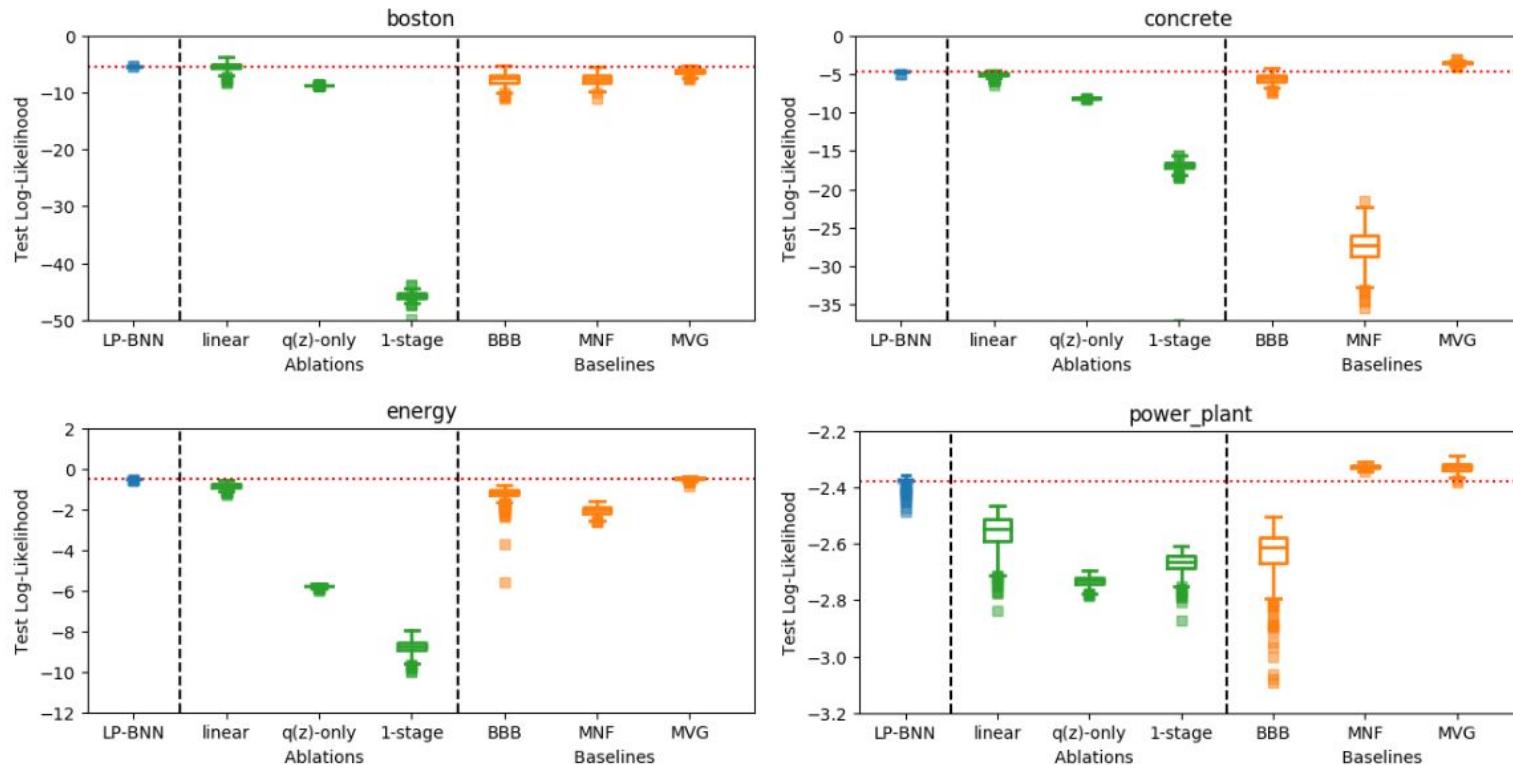
# Results: Generalization (Ablations)



1-stage			
linear			
$q(z)$ only			$q_{\lambda_z}(z)$

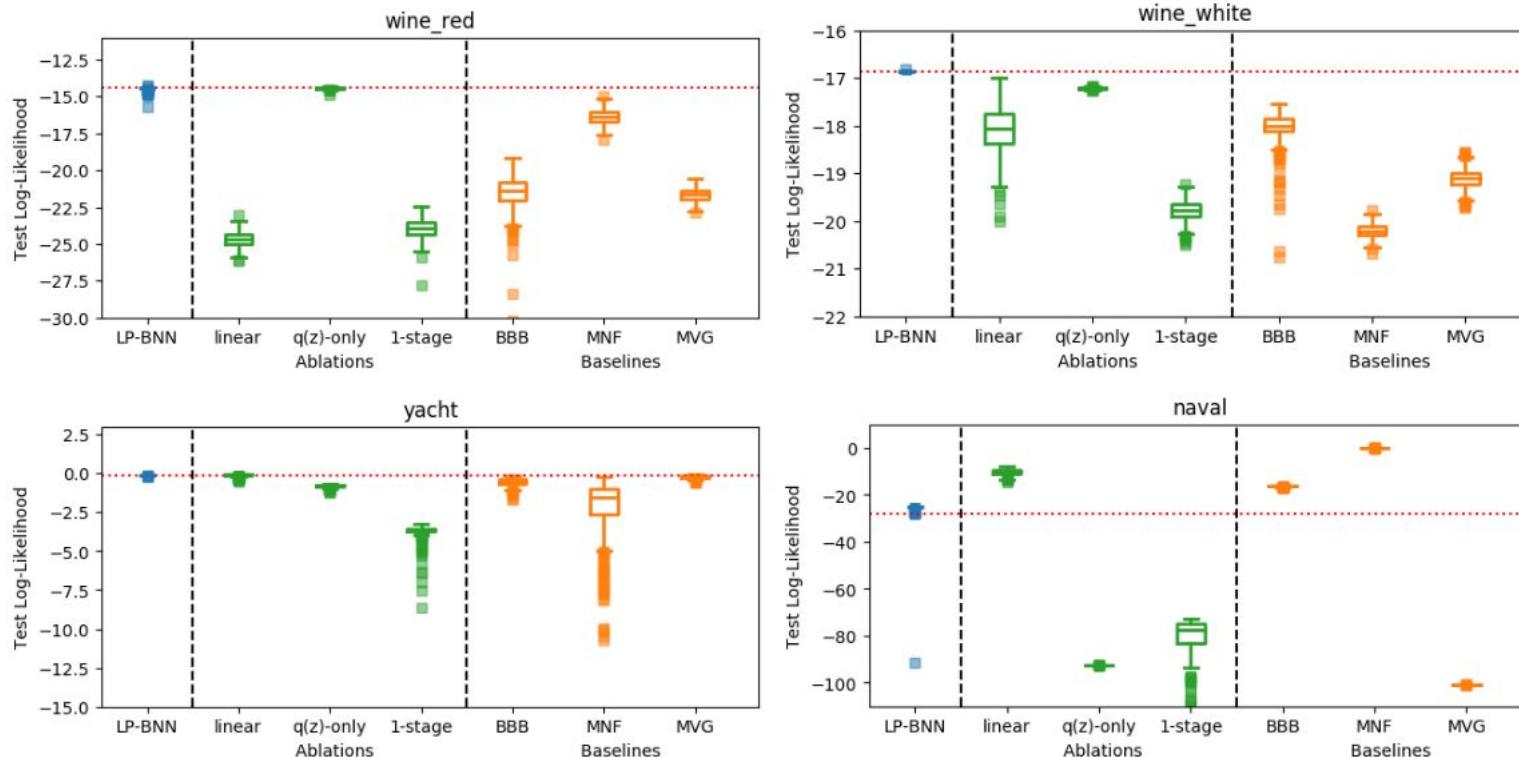
# Results: Generalization

<https://arxiv.org/abs/1811.07006>

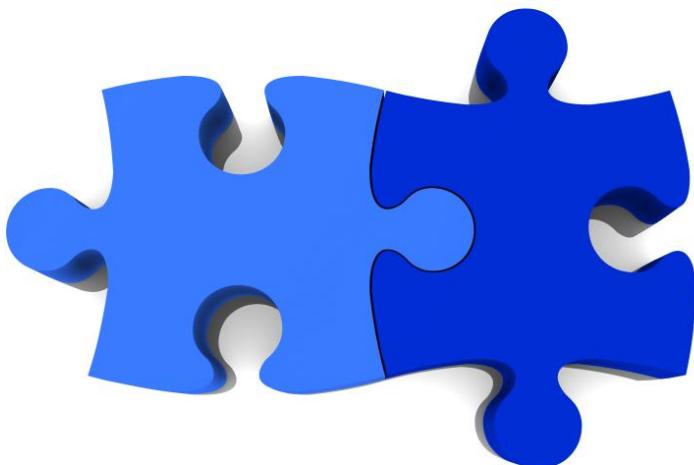


# Results: Generalization

<https://arxiv.org/abs/1811.07006>



# Conclusions



In this talk, two applications of LVMs

## 1. Data exploration

- a. Infinite latent feature model  
for heterogeneous datasets
- b. Global and group specific factors

<https://ivaleram.github.io/GLFM/>

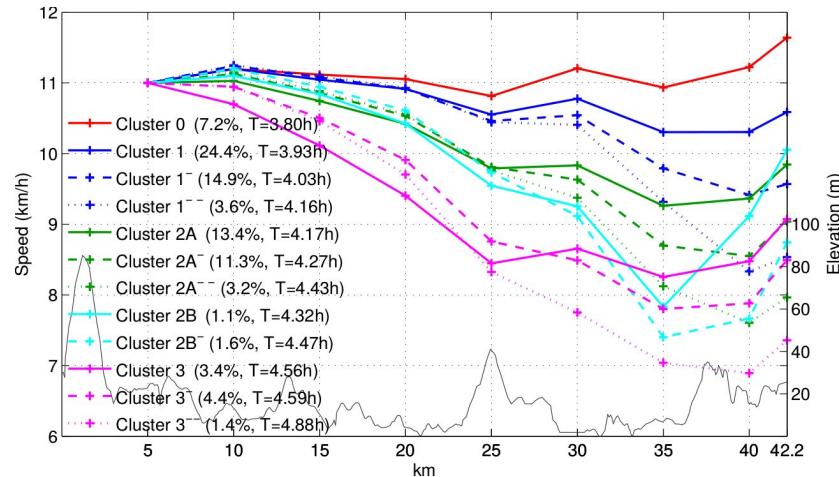
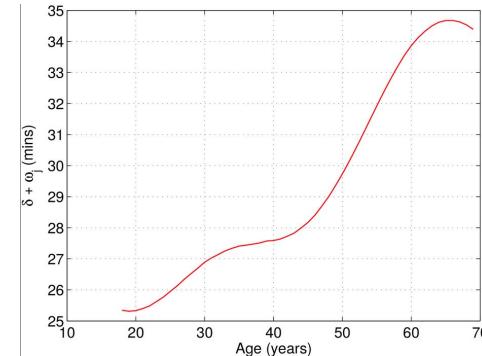
## 2. Uncertainty quantification

- a. Alternative modeling for BNNs
- b. Better approximate inference

<https://arxiv.org/abs/1811.07006>

# Other projects...

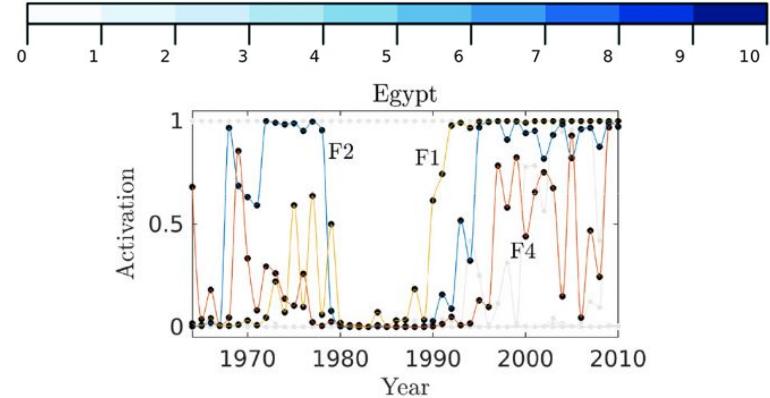
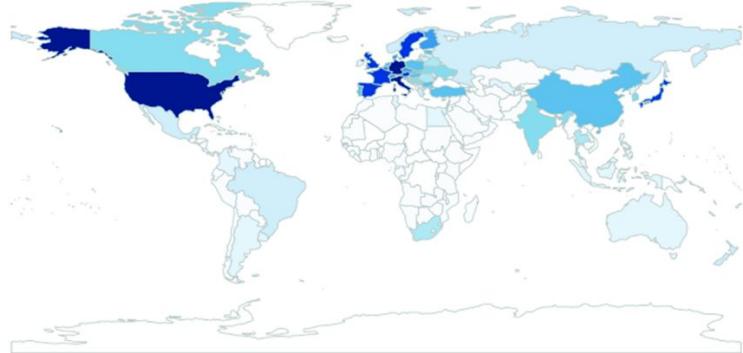
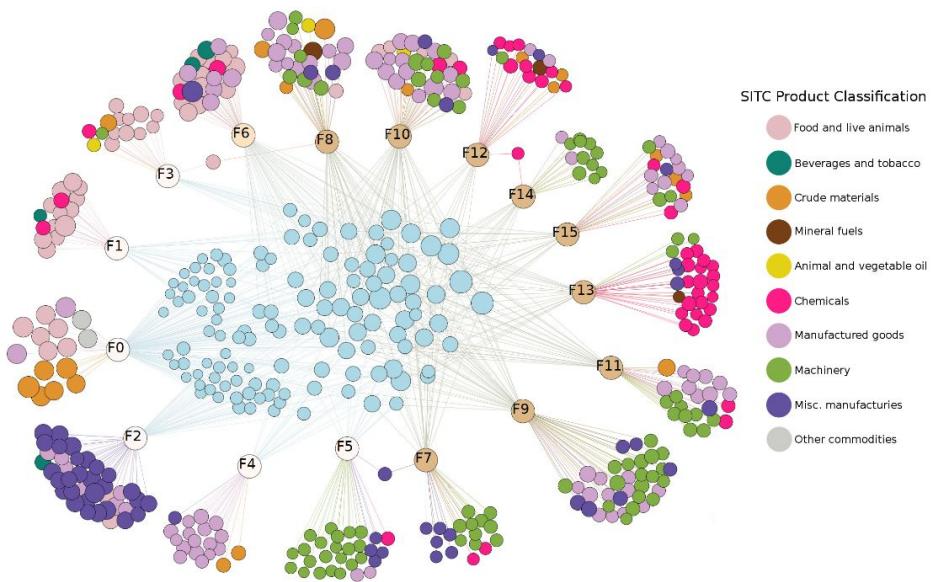
## Sport Science



M. F. Pradier, F. J. R. Ruiz, and F. Perez-Cruz. **Prior Design for Dependent Dirichlet Processes: An Application to Marathon Modeling**. *PlosONE*. 2016.

# Other projects...

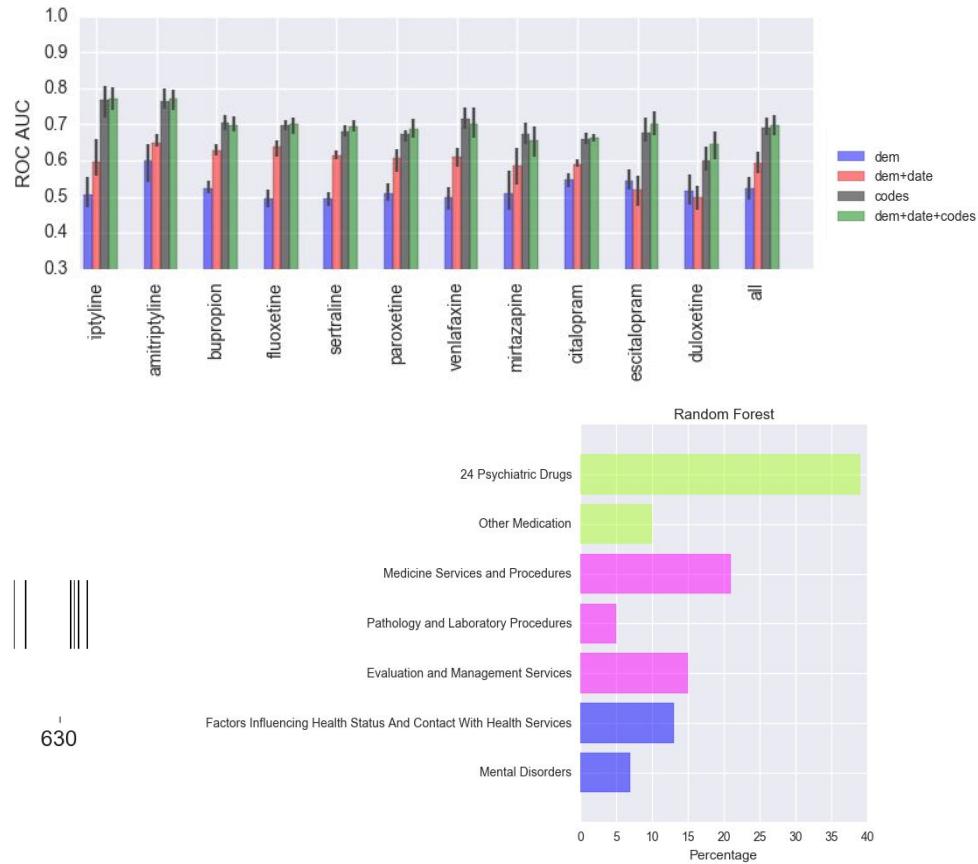
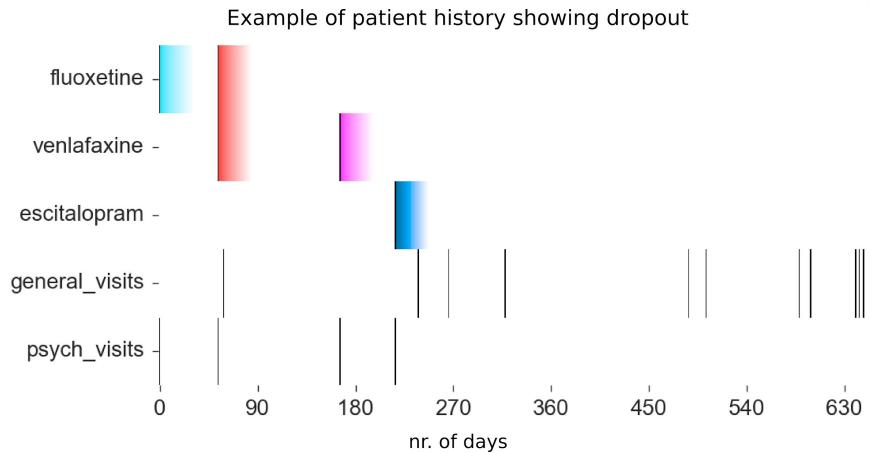
## Economics



Z. Utkovski, [M. F. Pradier](#), V. Stojkoski, L. Kocarev and F. Perez-Cruz. **Economic Complexity Unfolded: An Interpretable Model for the Productive Structure of Economies.** *PlosONE*. 2018.

# Other projects...

## Medicine: healthcare in psychiatry



M. F. Pradier, T. H. McCoy, M. Hughes, R. H. Perlis and F. Doshi-Velez. Predicting Treatment Discontinuation after Antidepressant Initiation. In submission to JAMA Psychiatry. 2018.

# Current research agenda



*Impact in real-world problems:*

- Personalize prescription of antidepressants

*ML research questions:*

- How to adapt model complexity automatically?
- How to better quantify model uncertainty?
- How to make models easy-to-interpret?
- How to combine expert knowledge with insights from data?

Contact: [melanie@seas.harvard.edu](mailto:melanie@seas.harvard.edu)

<https://melaniefp.github.io/>

## Special thanks to:

- Finale Doshi-Velez
- Weiwei Pan
- Michael Hughes
- All members of dtak!
- Francisco Rodriguez Ruiz
- Fernando Perez-Cruz
- Isabel Valera
- Maria Lomeli
- Zoubin Ghahramani
- Oscar Puig
- Francesca Milletti

# Thank you!



Weiwei Pan



Jiayu Yao



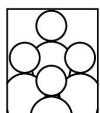
Soumya Ghosh



Maia Jacobs



Finale Doshi-Velez



Center for Research on  
Computation and Society

at Harvard John A. Paulson School of Engineering and Applied Sciences



HDSI

Harvard Data  
Science Initiative

<https://melaniefp.github.io/>

# Interpretable Machine Learning

## Interpretability

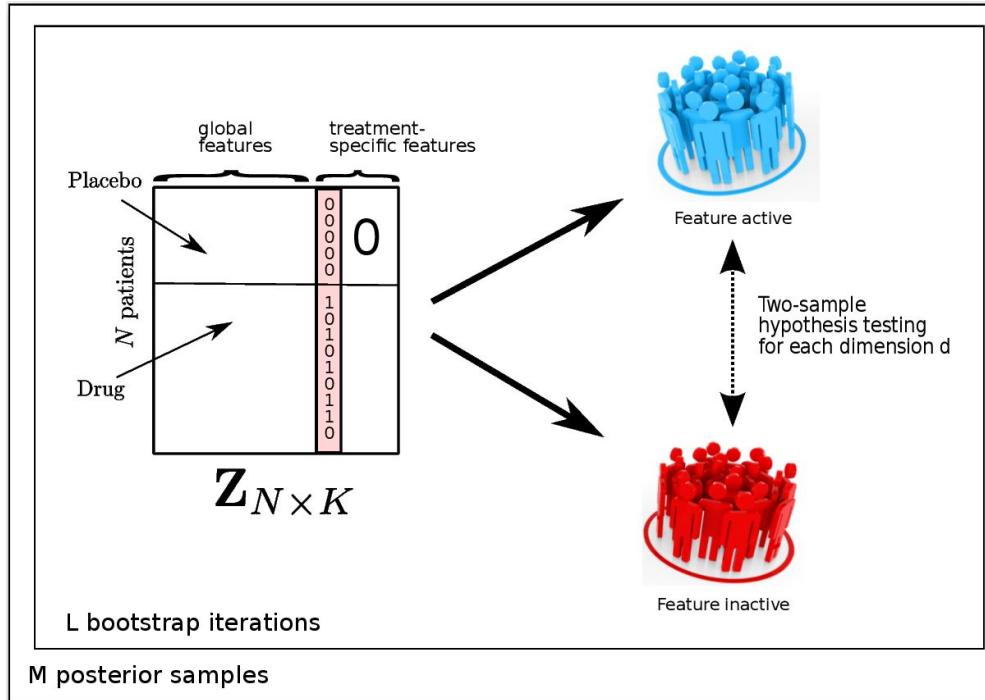
- “ability to explain or to present in understandable terms to a human” (Doshi-Velez and Kim, 2017)
- requirement in the 2018 EU General Data Protection Regulation (Goodman et.al. 2016)

## Interpretable Machine Learning

- Interpretable models to explain black-boxes
  - Local Interpretable Explanations (Ribeiro et.al, 2016)
  - Interpretable Decision Sets (Lakkaraju et.al, 2016)
- Interpretable models from scratch
  - Tree-regularization of deep models (Wu et.al, 2017)
  - Input-gradient regularization (Ross et.al, 2017)

In this talk, interpretability via prob. graphical models

# Statistical methodology for biomarker discovery



M. F. Pradier, B. Reis, L. Jukofsky, F. Milletti, T. Ohtomo, F. Perez-Cruz, and O. Puig. **Case-control Indian Buffet Process identifies biomarkers of response to Codrituzumab.** Accepted to *BMC Cancer*. 2019.

# Prediction-constrained Autoencoder

$$\begin{aligned} \{\boldsymbol{\theta}^*, \boldsymbol{\phi}^*\} = \operatorname{argmin}_{\boldsymbol{\theta}, \boldsymbol{\phi}} \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}) &= \min_{\boldsymbol{\theta}, \boldsymbol{\phi}} \left\{ \frac{1}{R} \sum_{r=1}^R \left( \mathbf{w_c}^{(r)} - g_{\boldsymbol{\phi}} \left( f_{\boldsymbol{\theta}} \left( \mathbf{w_c}^{(r)} \right) \right) + \gamma^{(r)} \right)^2 \right. \\ &\quad \left. + \beta \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[ \frac{1}{R} \sum_{r=1}^R \log p(y|x, g_{\boldsymbol{\phi}} \left( f_{\boldsymbol{\theta}} \left( \mathbf{w_c}^{(r)} \right) \right)) \right] \right\}, \end{aligned}$$