

ARE YOU EXPLOITING YOUR ASSUMPTIONS?
TOWARDS EXPRESSIVE PRIORS FOR BIOMARKER
DISCOVERY AND FUNCTIONAL PREDICTION

Melanie F. Pradier

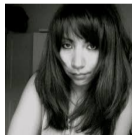
Harvard University

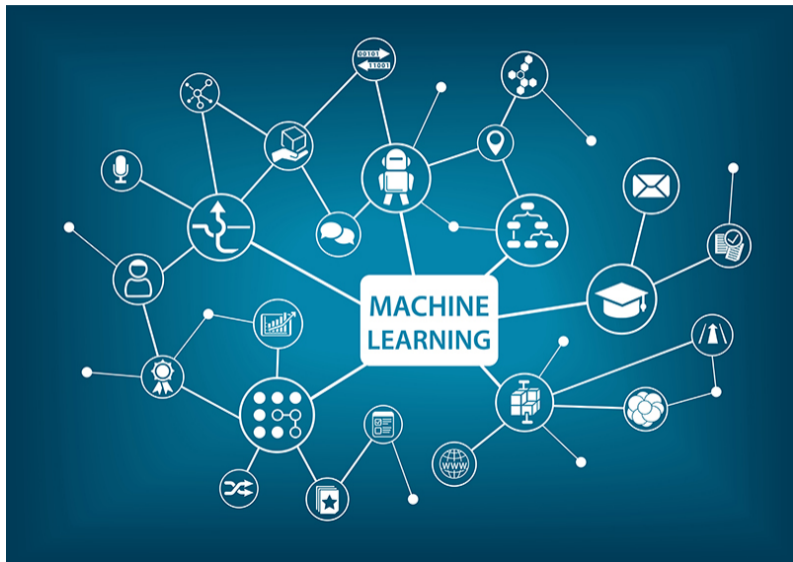
February 18th, 2020

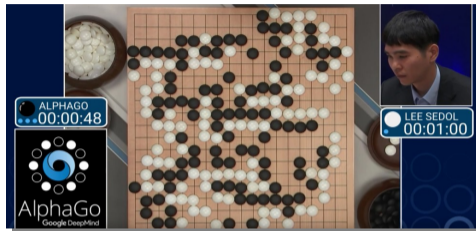
ARE YOU EXPLOITING YOUR ASSUMPTIONS?

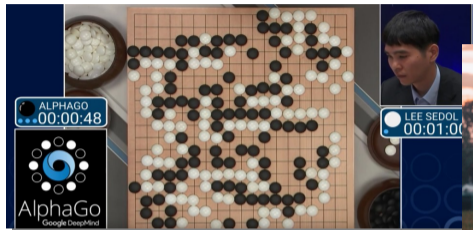
TOWARDS EXPRESSIVE PRIORS FOR BIOMARKER DISCOVERY AND FUNCTIONAL PREDICTION

Research in collaboration
with...

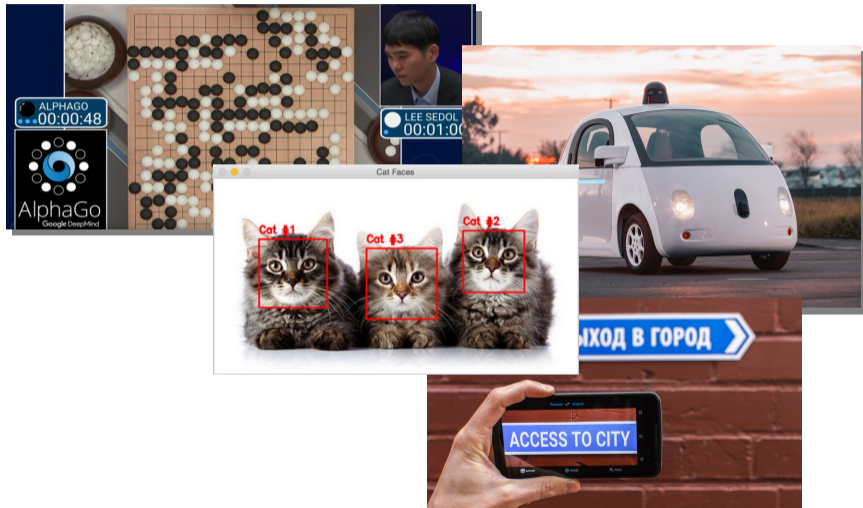






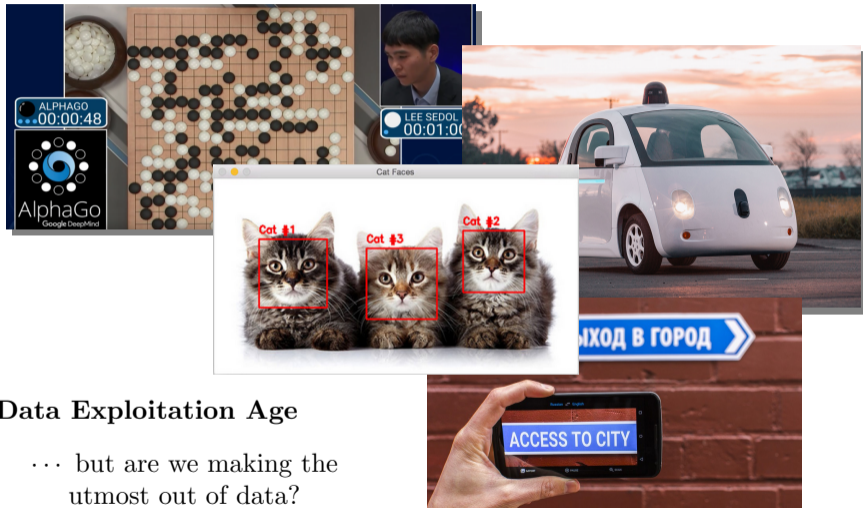








Data Exploitation Age



Data Exploitation Age

... but are we making the utmost out of data?

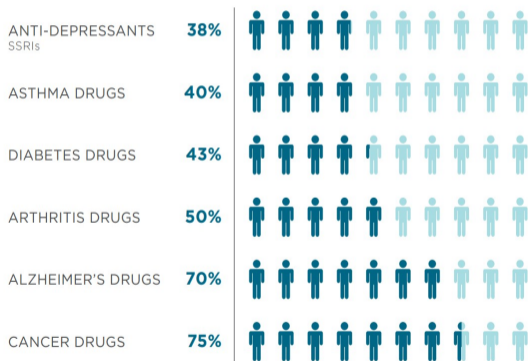
ARE WE MAKING THE UTMOST OUT OF DATA?

AN EXAMPLE: PERSONALIZED MEDICINE

ARE WE MAKING THE UTMOST OUT OF DATA?

AN EXAMPLE: PERSONALIZED MEDICINE

Percentage of the patient population for which a particular drug in a class is ineffective, on average



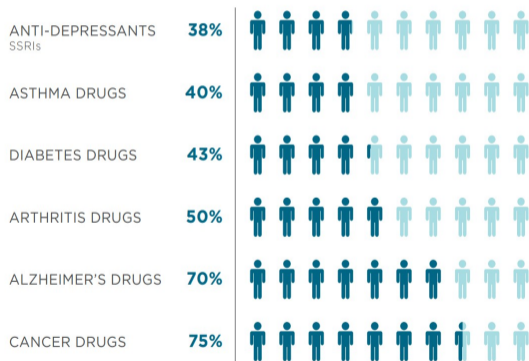
Source: Brian B. Spear, Margo Heath-Chiozzi, Jeffrey Huff, "Clinical Trends in Molecular Medicine," Volume 7, Issue 5, 1 May 2001, pages 201-204.

ARE WE MAKING THE UTMOST OUT OF DATA?

AN EXAMPLE: PERSONALIZED MEDICINE



Percentage of the patient population for which a particular drug in a class is ineffective, on average



Source: Brian B. Spear, Margo Heath-Chiozzi, Jeffrey Huff, "Clinical Trends in Molecular Medicine," Volume 7, Issue 5, 1 May 2001, pages 201-204.

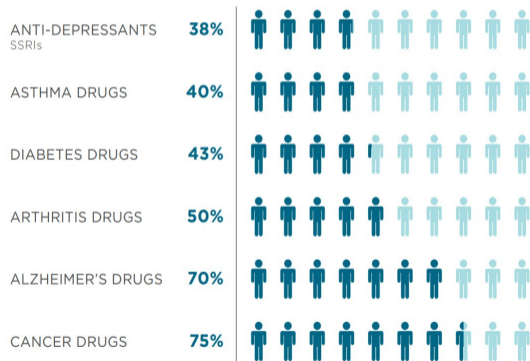
ARE WE MAKING THE UTMOST OUT OF DATA?

AN EXAMPLE: PERSONALIZED MEDICINE



CHALLENGES

Percentage of the patient population for which a particular drug in a class is ineffective, on average



► Complexity

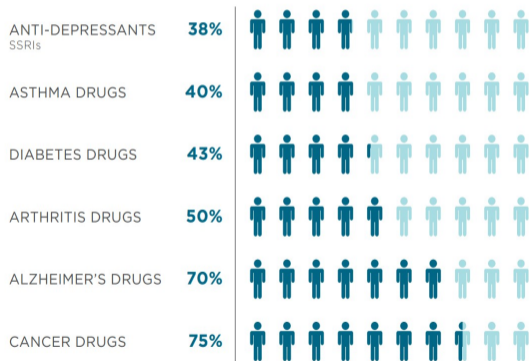
Source: Brian B. Spear, Margo Heath-Chiozzi, Jeffrey Huff, "Clinical Trends in Molecular Medicine," Volume 7, Issue 5, 1 May 2001, pages 201-204.

ARE WE MAKING THE UTMOST OUT OF DATA?

AN EXAMPLE: PERSONALIZED MEDICINE



Percentage of the patient population for which a particular drug in a class is ineffective, on average



CHALLENGES

- ▶ Complexity
- ▶ Noise, missing data

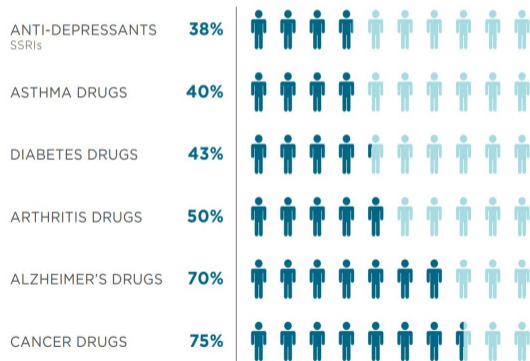
Source: Brian B. Spear, Margo Heath-Chiozzi, Jeffrey Huff, "Clinical Trends in Molecular Medicine," Volume 7, Issue 5, 1 May 2001, pages 201-204.

ARE WE MAKING THE UTMOST OUT OF DATA?

AN EXAMPLE: PERSONALIZED MEDICINE



Percentage of the patient population for which a particular drug in a class is ineffective, on average



CHALLENGES

- ▶ Complexity
- ▶ Noise, missing data
- ▶ *Small data within big data*
- ▶ ...

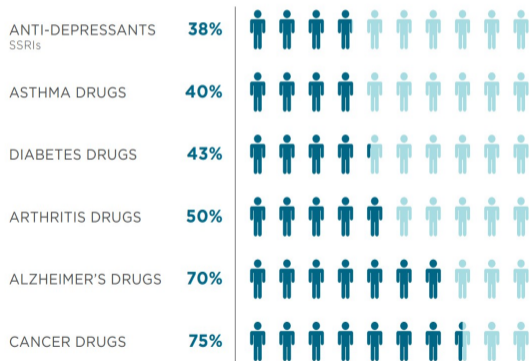
Source: Brian B. Spear, Margo Heath-Chiozzi, Jeffrey Huff, "Clinical Trends in Molecular Medicine," Volume 7, Issue 5, 1 May 2001, pages 201-204.

ARE WE MAKING THE UTMOST OUT OF DATA?

AN EXAMPLE: PERSONALIZED MEDICINE



Percentage of the patient population for which a particular drug in a class is ineffective, on average



CHALLENGES

- ▶ Complexity
- ▶ Noise, missing data
- ▶ *Small data within big data*
- ▶ ...
- ▶ **Need to understand**
→ data exploration

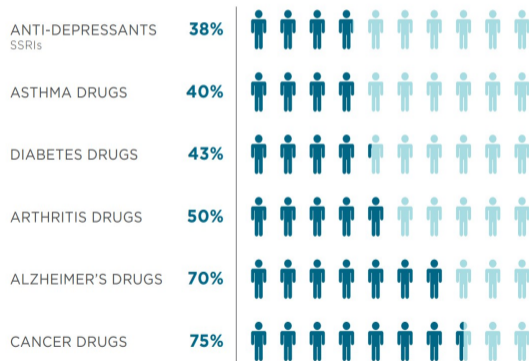
Source: Brian B. Spear, Margo Heath-Chiozzi, Jeffrey Huff, "Clinical Trends in Molecular Medicine," Volume 7, Issue 5, 1 May 2001, pages 201-204.

ARE WE MAKING THE UTMOST OUT OF DATA?

AN EXAMPLE: PERSONALIZED MEDICINE



Percentage of the patient population for which a particular drug in a class is ineffective, on average



CHALLENGES

- ▶ Complexity
- ▶ Noise, missing data
- ▶ *Small data within big data*
- ▶ ...
- ▶ **Need to understand**
→ data exploration

How can ML systems help experts “understand” data?

Source: Brian B. Spear, Margo Heath-Chiozzi, Jeffrey Huff, "Clinical Trends in Molecular Medicine," Volume 7, Issue 5, 1 May 2001, pages 201-204.

EXPRESSIVITY-INTERPRETABILITY LOOP

INTERPRETABILITY

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

EXPRESSIVITY-INTERPRETABILITY LOOP

INTERPRETABILITY

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

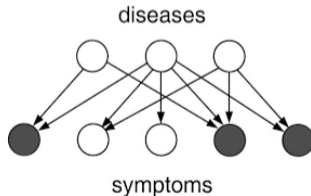
EXPRESSIVITY

- ▶ ability to encode assumptions/desiderata into the model

In this talk, expressivity/interpretability via probabilistic graphical models

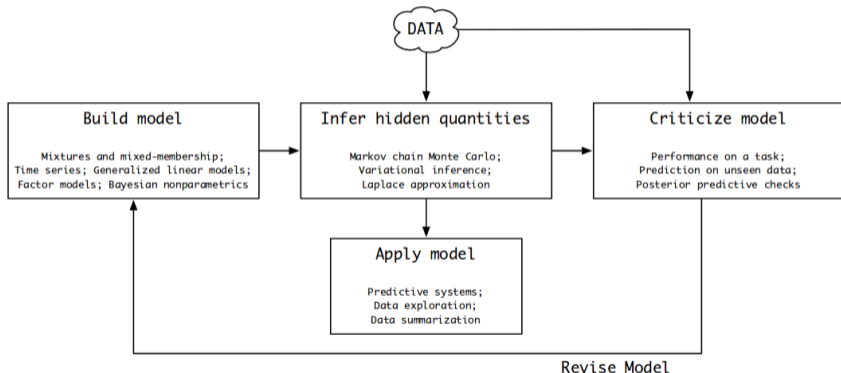
WHY PROBABILISTIC GRAPHICAL MODELS?

- ▶ Generative model \equiv unsupervised approach, goal is to model $p(\text{Data})$
- ▶ Graphical model for multidisciplinary research
- ▶ Assumptions and desiderata through *latent variables*



WHY PROBABILISTIC GRAPHICAL MODELS?

THE “BOX’S LOOP” (BLEI, 2014)



OUTLINE

- ▶ Overview
- ▶ Goal I: Biomarker discovery
- ▶ Goal II: Functional prediction
- ▶ Wrap-up

OUR FOCUS: BIOMARKER DISCOVERY

DEF: "ANY VARIABLE THAT CAN BE USED AS AN INDICATOR OF A PARTICULAR DISEASE STATE".

OUR FOCUS: BIOMARKER DISCOVERY

DEF: "ANY VARIABLE THAT CAN BE USED AS AN INDICATOR OF A PARTICULAR DISEASE STATE".

Biomarkers are used everywhere!!

OUR FOCUS: BIOMARKER DISCOVERY

DEF: "ANY VARIABLE THAT CAN BE USED AS AN INDICATOR OF A PARTICULAR DISEASE STATE".

Biomarkers are used everywhere!!

SOME EXAMPLES

- ▶ 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

OUR FOCUS: BIOMARKER DISCOVERY

DEF: "ANY VARIABLE THAT CAN BE USED AS AN INDICATOR OF A PARTICULAR DISEASE STATE".

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

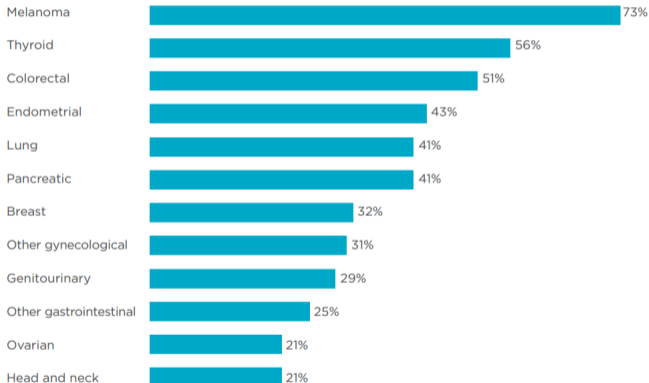
What's the
optimal dose
for my body?

Recurrence

Will the
cancer
return?

BIOMARKERS AS POTENTIAL TARGETS FOR NEW DRUGS

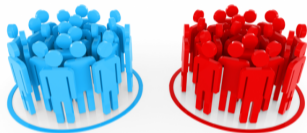
TACKLING TUMORS: Percentage of patients whose tumors were driven by certain genetic mutations that could be targets for specific drugs, by types of cancer.



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

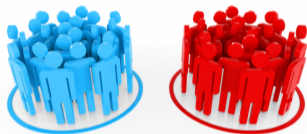
PROBLEM FORMULATION

BIOMARKER DISCOVERY IN CLINICAL TRIALS



PROBLEM FORMULATION

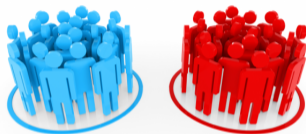
BIOMARKER DISCOVERY IN CLINICAL TRIALS



We want to discover:

PROBLEM FORMULATION

BIOMARKER DISCOVERY IN CLINICAL TRIALS



We want to discover:

1. Indicators of disease progression: prognostic biomarkers

PROBLEM FORMULATION

BIOMARKER DISCOVERY IN CLINICAL TRIALS

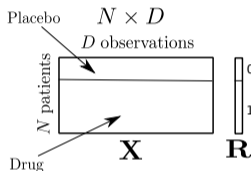
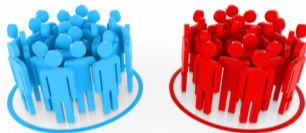


We want to discover:

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

PROBLEM FORMULATION

BIOMARKER DISCOVERY IN CLINICAL TRIALS



We want to discover:

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

APPLICATION: IMMUNOTHERAPY FOR LIVER CANCER

[ABOU-ALFA ET.AL, 2016]

SO FAR...

- ▶ No evidence for treatment effectiveness
- ▶ Hypothesis: drug exposure as confounder

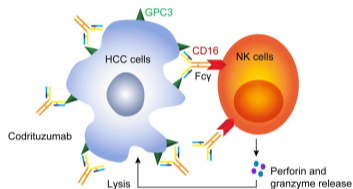
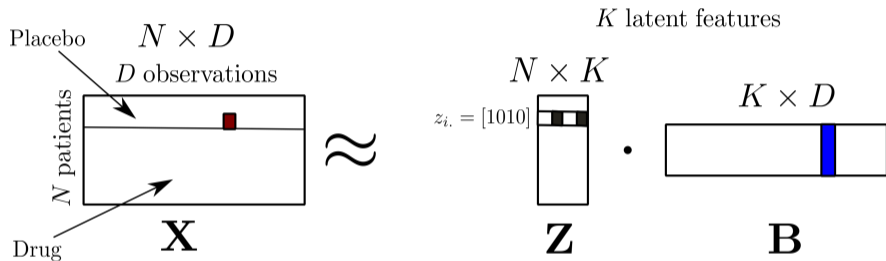


Diagram: Mechanism of codrituzumab-induced antibody-dependent cytotoxicity through the interaction of Fc CD16 in NK cells

HOW TO DEAL WITH...

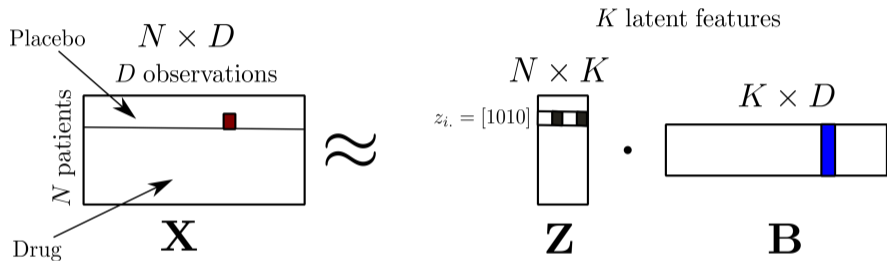
- ▶ data complexity?
- ▶ data heterogeneity?
- ▶ drug effect vs natural response?

OUR APPROACH: LATENT FEATURE MODEL



■ $x_{id} = 173 \text{ ml/dL} = 73 + 0 + 100 \text{ ml/dL}$

OUR APPROACH: LATENT FEATURE MODEL



$$\blacksquare x_{id} = 173 \text{ ml/dL} = 73 + 0 + 100 \text{ ml/dL}$$

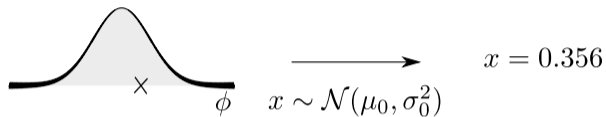
Note: Correlation does not imply causality!

HOW TO DEAL WITH COMPLEXITY AND PATIENT HETEROGENEITY?

INDIAN BUFFET PROCESS [GHAHRAMANI ET.AL, 2006]

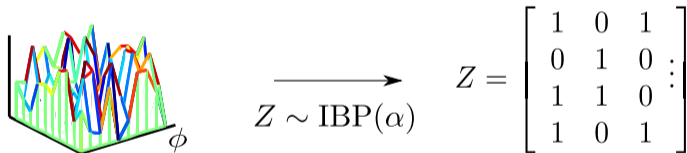
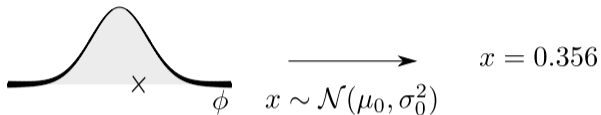
HOW TO DEAL WITH COMPLEXITY AND PATIENT HETEROGENEITY?

INDIAN BUFFET PROCESS [GHAHRAMANI ET.AL, 2006]



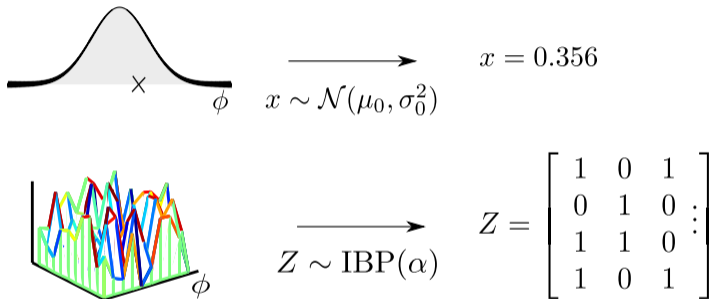
HOW TO DEAL WITH COMPLEXITY AND PATIENT HETEROGENEITY?

INDIAN BUFFET PROCESS [GHAHRAMANI ET.AL, 2006]



HOW TO DEAL WITH COMPLEXITY AND PATIENT HETEROGENEITY?

INDIAN BUFFET PROCESS [GHAHRAMANI ET.AL, 2006]



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

HOW TO DEAL WITH COMPLEXITY AND PATIENT HETEROGENEITY?

INDIAN BUFFET PROCESS [GHAHRAMANI ET.AL, 2006]



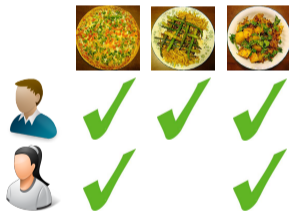
HOW TO DEAL WITH COMPLEXITY AND PATIENT HETEROGENEITY?

INDIAN BUFFET PROCESS [GHAHRAMANI ET.AL, 2006]



HOW TO DEAL WITH COMPLEXITY AND PATIENT HETEROGENEITY?

INDIAN BUFFET PROCESS [GHAHRAMANI ET.AL, 2006]



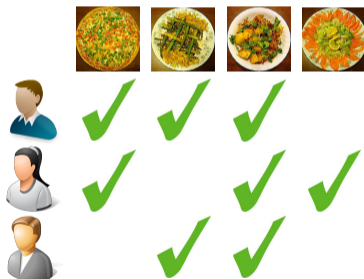
HOW TO DEAL WITH COMPLEXITY AND PATIENT HETEROGENEITY?

INDIAN BUFFET PROCESS [GHAHRAMANI ET.AL, 2006]



HOW TO DEAL WITH COMPLEXITY AND PATIENT HETEROGENEITY?

INDIAN BUFFET PROCESS [GHAHRAMANI ET.AL, 2006]



HOW TO DEAL WITH COMPLEXITY AND PATIENT HETEROGENEITY?

INDIAN BUFFET PROCESS [GHAHRAMANI ET.AL, 2006]



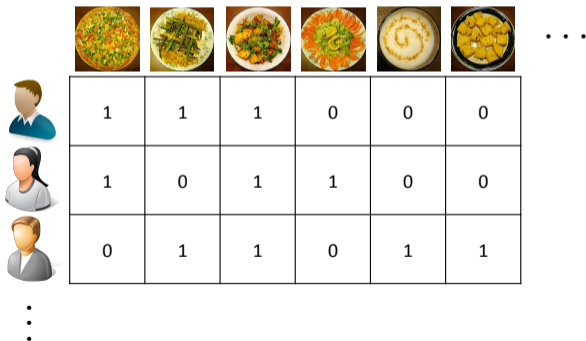
HOW TO DEAL WITH COMPLEXITY AND PATIENT HETEROGENEITY?

INDIAN BUFFET PROCESS [GHAHRAMANI ET.AL, 2006]












HOW TO DEAL WITH COMPLEXITY AND PATIENT HETEROGENEITY?

INDIAN BUFFET PROCESS [GHAHRAMANI ET.AL, 2006]

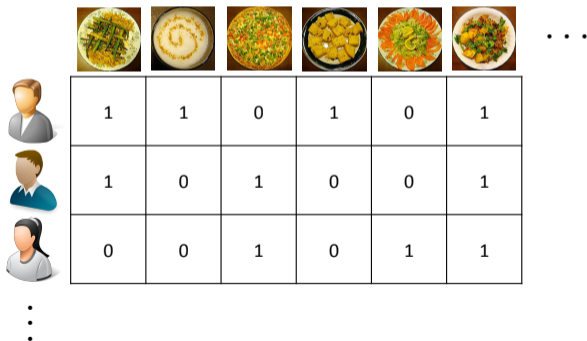


The diagram illustrates the Indian Buffet Process. It features a grid where rows represent patients and columns represent dishes. Each cell contains a binary value (0 or 1) indicating whether a patient has sampled a dish. Above the grid are six images of different Indian dishes: a vegetable salad, a vegetable stir-fry, a vegetable and chickpea salad, a vegetable and carrot salad, a rice dish with a yellow sauce, and a vegetable and chickpea salad. To the right of the grid are three dots indicating more dishes. To the left of the grid are three icons representing different patients, with three dots below them indicating more patients.










							...
1	1	1	0	0	0		
	1	0	1	1	0	0	
	0	1	1	0	1	1	
...							

HOW TO DEAL WITH COMPLEXITY AND PATIENT HETEROGENEITY?

INDIAN BUFFET PROCESS [GHAHRAMANI ET.AL, 2006]

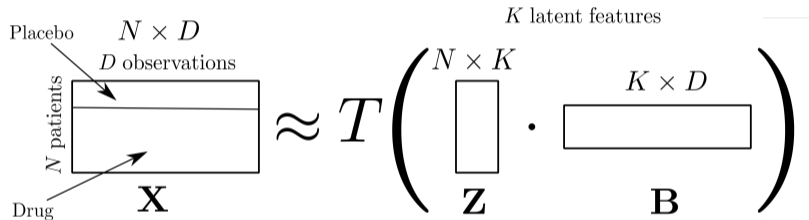


The diagram illustrates the Indian Buffet Process. At the top, six different Indian dishes are shown in a row, followed by an ellipsis. Below the dishes is a grid where each row represents a patient and each column represents a dish. The grid contains binary values (0 or 1) indicating whether a patient has selected a particular dish. To the left of the grid are three patient icons, and below them is a vertical ellipsis. To the right of the grid is a horizontal ellipsis.

							...
1	1	0	1	0	1		
	1	0	1	0	0	1	
	0	0	1	0	1	1	
...							

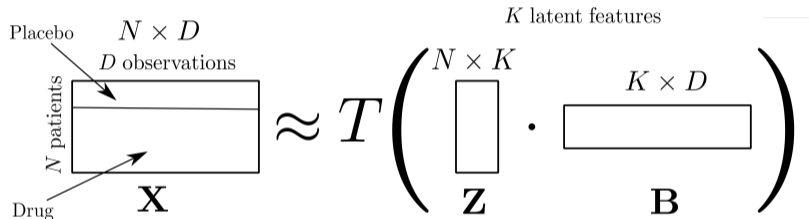
HOW ABOUT FEATURE HETEROGENEITY?

GENERAL LATENT FEATURE MODEL (GLFM) [VALERA ET.AL, 2020]



HOW ABOUT FEATURE HETEROGENEITY?

GENERAL LATENT FEATURE MODEL (GLFM) [VALERA ET.AL, 2020]

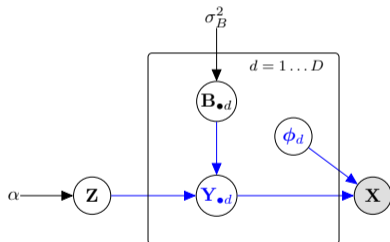


- ▶ K potentially unbounded
- ▶ Link functions T_d for each feature d

HOW ABOUT FEATURE HETEROGENEITY?

GENERAL LATENT FEATURE MODEL (GLFM) [VALERA ET.AL, 2020]

Latent feature model for heterogeneous datasets



- ▶ Link functions T_d depend on type of data for each feature d

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

$$y_{nd} | \mathbf{Z}, \mathbf{B} \sim \mathcal{N}(\mathbf{Z}_{n\bullet} \mathbf{B}_{\bullet,d}, \sigma_y^2)$$

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

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

GLFM PACKAGE

- ▶ Open-source python/matlab/R code
- ▶ Deals with Heterogeneous datasets

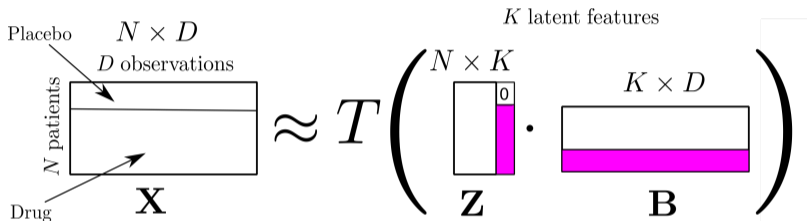
<https://github.com/ivaleraM/GLFM>

HOW TO DISTINGUISH DRUG EFFECT VS NATURAL RESPONSE?

CASE-CONTROL INDIAN BUFFET PROCESS (C-IBP) [PRADIER ET.AL, 2019]

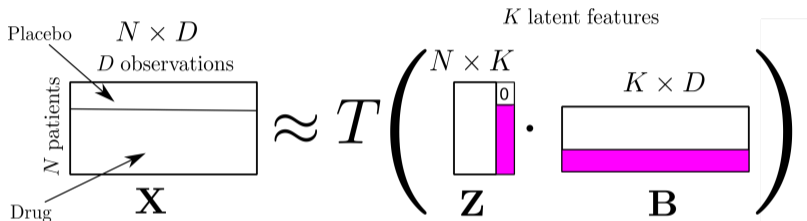
HOW TO DISTINGUISH DRUG EFFECT VS NATURAL RESPONSE?

CASE-CONTROL INDIAN BUFFET PROCESS (C-IBP) [PRADIER ET.AL, 2019]



HOW TO DISTINGUISH DRUG EFFECT VS NATURAL RESPONSE?

CASE-CONTROL INDIAN BUFFET PROCESS (C-IBP) [PRADIER ET.AL, 2019]

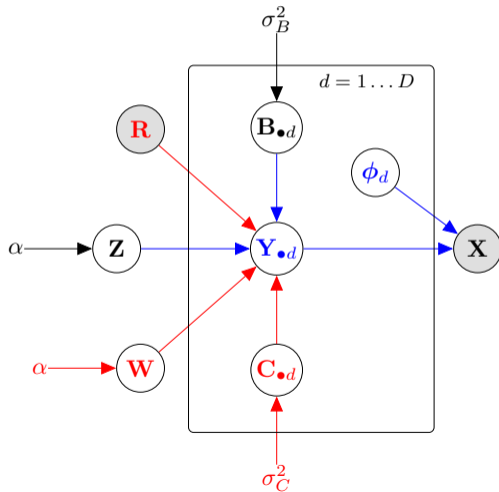


TREATMENT-SPECIFIC LATENT FEATURES

- ▶ can only activate for patients in treatment arm
- ▶ number learned automatically

HOW TO DISTINGUISH DRUG EFFECT VS NATURAL RESPONSE?

CASE-CONTROL INDIAN BUFFET PROCESS (C-IBP) [PRADIER ET.AL, 2019]



R_n : drug indicator per patient n

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

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

$$\mathcal{N}(\mathbf{Z}_n \bullet \mathbf{B}_{\bullet d} + \mathbb{1}[\mathbf{R}_n = 1] \mathbf{W}_n \bullet \mathbf{C}_{\bullet d}, \sigma_y^2)$$

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

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

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

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

- ▶ **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

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

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

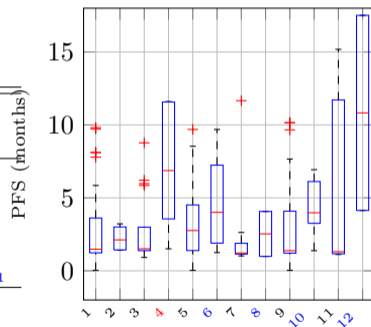
Sub-population	Drug Identifier				Size (number of patients)	Mean PFS (months)	Median PFS (months)
	F1	F2	F3				
1.	0	0	0	0	33.37	3.06	1.65
2.	0	0	1	0	4.07	2.29	2.24
3.	0	1	0	0	17.84	2.72	1.81
4.	0	1	1	0	4.72	7.05	7.18
5.	1	0	0	0	51.52	3.22	2.55
6.	1	0	0	1	16.77	4.17	3.65
7.	1	0	1	0	8.38	1.74	1.33
8.	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

RESULTS: SUBPOPULATIONS

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

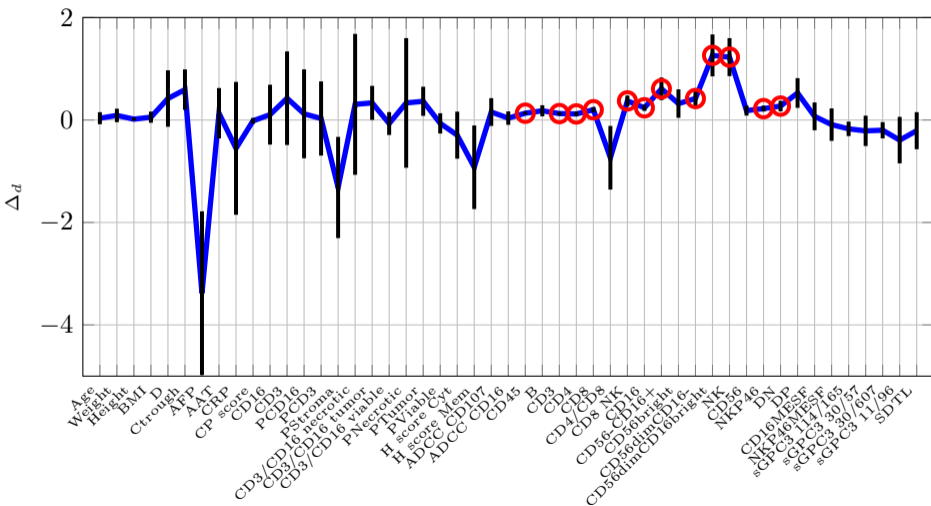
- ▶ 180 patients: 60 took a placebo, 120 took the drug
- ▶ PFS: Progression Free Survival

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



RESULTS: BIOMARKER DISCOVERY

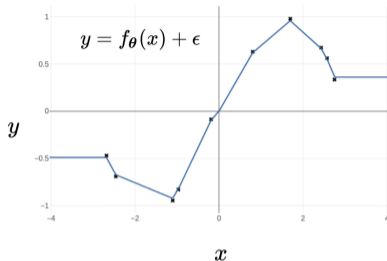
TREATMENT-SPECIFIC FEATURE F3



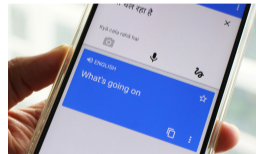
OUTLINE

- ▶ Overview
- ▶ Goal I: Biomarker discovery
- ▶ Goal II: Functional prediction
- ▶ Wrap-up

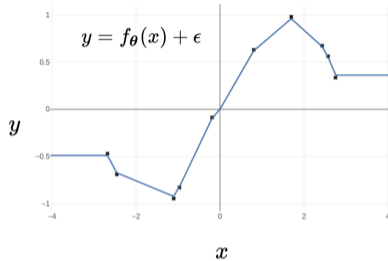
NEURAL NETWORKS (NNs) AS UNIVERSAL APPROXIMATORS



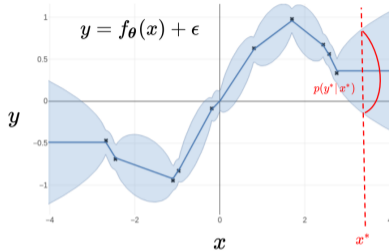
Several success stories...



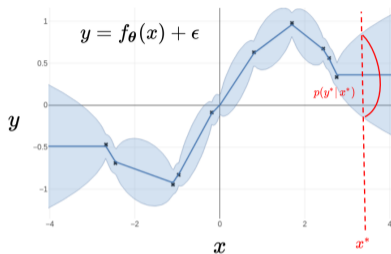
BUT WHAT IF STAKES ARE HIGH?



BUT WHAT IF STAKES ARE HIGH?



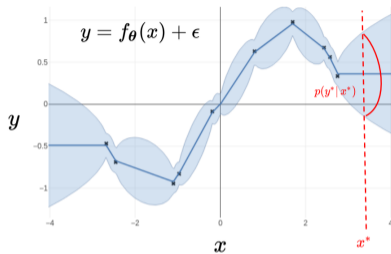
BUT WHAT IF STAKES ARE HIGH?



Uncertainty estimation becomes crucial!



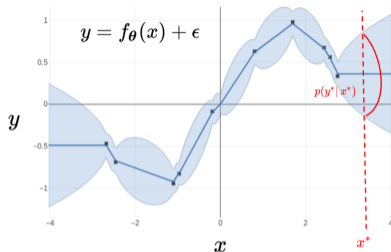
SOMETIMES WE HAVE A PRIORI FUNCTIONAL KNOWLEDGE...



Some examples of assumptions:

- ▶ Range of heart rate at rest between 60-100 bpm.
- ▶ Slow/fast variation of air pollutant
- ▶ Volatility of stock market

SOMETIMES WE HAVE A PRIORI FUNCTIONAL KNOWLEDGE...



Some examples of assumptions:

- ▶ Range of heart rate at rest between 60-100 bpm.
- ▶ Slow/fast variation of air pollutant
- ▶ Volatility of stock market

How can we incorporate such functional desiderata into the model?

AN EASY WAY TO SPECIFY FUNCTIONAL DESIDERATA: GAUSSIAN PROCESSES (GPs)

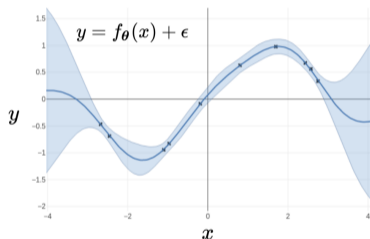
Definition: a Gaussian process is a collection of random variables, any finite number of which have (consistent) Gaussian distributions.

$$f \sim \mathcal{N}(\mu(\cdot), k(\cdot, \cdot))$$

Example: RBF kernel as covariance function:

$$k(x, x') = \sigma^2 \exp\left(-\frac{(x - x')^2}{2\gamma^2}\right)$$

- ▶ Stationarity
- ▶ Lengthscale
- ▶ Amplitude variance



AN EASY WAY TO SPECIFY FUNCTIONAL DESIDERATA: GAUSSIAN PROCESSES (GPs)

Definition: a Gaussian process is a collection of random variables, any finite number of which have (consistent) Gaussian distributions.

$$f \sim \mathcal{N}(\mu(\cdot), k(\cdot, \cdot))$$

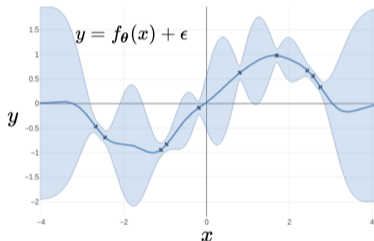
Example: RBF kernel as covariance function:

$$k(x, x') = \sigma^2 \exp\left(-\frac{(x - x')^2}{2\gamma^2}\right)$$

► Stationarity

► Lengthscale

► Amplitude variance



AN EASY WAY TO SPECIFY FUNCTIONAL DESIDERATA: GAUSSIAN PROCESSES (GPs)

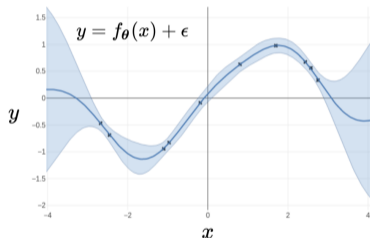
Definition: a Gaussian process is a collection of random variables, any finite number of which have (consistent) Gaussian distributions.

$$f \sim \mathcal{N}(\mu(\cdot), k(\cdot, \cdot))$$

Example: RBF kernel as covariance function:

$$k(x, x') = \sigma^2 \exp\left(-\frac{(x - x')^2}{2\gamma^2}\right)$$

▶ Stationarity



▶ Lengthscale

▶ Amplitude variance

AN EASY WAY TO SPECIFY FUNCTIONAL DESIDERATA: GAUSSIAN PROCESSES (GPs)

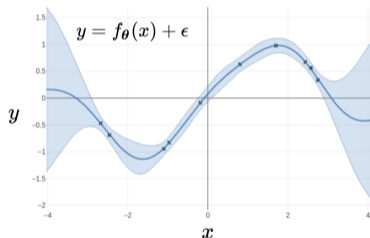
Definition: a Gaussian process is a collection of random variables, any finite number of which have (consistent) Gaussian distributions.

$$f \sim \mathcal{N}(\mu(\cdot), k(\cdot, \cdot))$$

Example: RBF kernel as covariance function:

$$k(x, x') = \sigma^2 \exp\left(-\frac{(x - x')^2}{2\gamma^2}\right)$$

- ▶ Stationarity
- ▶ Lengthscale
- ▶ Amplitude variance



AN EASY WAY TO SPECIFY FUNCTIONAL DESIDERATA: GAUSSIAN PROCESSES (GPs)

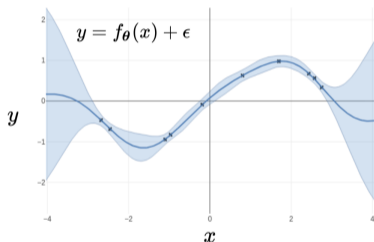
Definition: a Gaussian process is a collection of random variables, any finite number of which have (consistent) Gaussian distributions.

$$f \sim \mathcal{N}(\mu(\cdot), k(\cdot, \cdot))$$

Example: RBF kernel as covariance function:

$$k(x, x') = \sigma^2 \exp\left(-\frac{(x - x')^2}{2\gamma^2}\right)$$

- ▶ Stationarity
- ▶ Lengthscale
- ▶ Amplitude variance



GPs ARE GREAT, BUT WHAT IF I STILL WANT A NN?

Benefits of NN approaches:

- ▶ widely used (many tools available)
- ▶ parametric expression
- ▶ fast at evaluation time

GPs ARE GREAT, BUT WHAT IF I STILL WANT A NN?

Benefits of NN approaches:

- ▶ widely used (many tools available)
- ▶ parametric expression
- ▶ fast at evaluation time

KEY RESEARCH QUESTIONS:

1. Can we design Bayesian NN priors that encode **stationarity properties** like a GP while retaining the benefits of neural networks?
2. Can we easily specify lengthscale and amplitude variance in a **decoupled** fashion?

BAYESIAN NEURAL NETWORKS

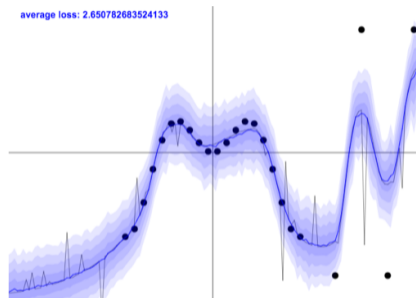
- ▶ Assume prior on network parameters
- ▶ Most common, i.i.d Gaussians

$$\mathbf{y} = f_{\boldsymbol{\theta}}(\mathbf{x}) + \boldsymbol{\epsilon}$$

$$\boldsymbol{\epsilon} \sim \mathcal{N}(0, \sigma_y^2 I)$$

$$\theta_i \sim \mathcal{N}(0, \sigma_{\theta}^2 I) \quad \forall i$$

- ▶ $p(\boldsymbol{\theta}) \implies p(f)$



(Yarin Gal blog)

BAYESIAN NEURAL NETWORKS

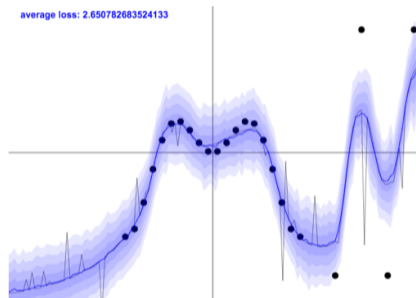
- ▶ Assume prior on network parameters
- ▶ Most common, i.i.d Gaussians

$$\mathbf{y} = f_{\boldsymbol{\theta}}(\mathbf{x}) + \epsilon$$

$$\epsilon \sim \mathcal{N}(0, \sigma_y^2 I)$$

$$\theta_i \sim \mathcal{N}(0, \sigma_{\theta}^2 I) \quad \forall i$$

- ▶ $p(\boldsymbol{\theta}) \implies p(f)$



(Yarin Gal blog)

- ▶ But what does a prior over weights mean in function space?

BAYESIAN NEURAL NETWORKS

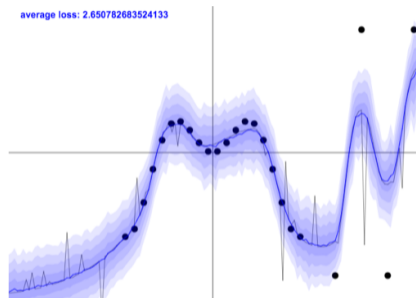
- ▶ Assume prior on network parameters
- ▶ Most common, i.i.d Gaussians

$$\mathbf{y} = f_{\boldsymbol{\theta}}(\mathbf{x}) + \epsilon$$

$$\epsilon \sim \mathcal{N}(0, \sigma_y^2 I)$$

$$\theta_i \sim \mathcal{N}(0, \sigma_{\theta}^2 I) \quad \forall i$$

- ▶ $p(\boldsymbol{\theta}) \implies p(f)$

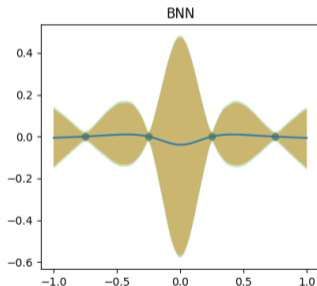


(Yarin Gal blog)

- ▶ But what does a prior over weights mean in function space?
Hard to know!

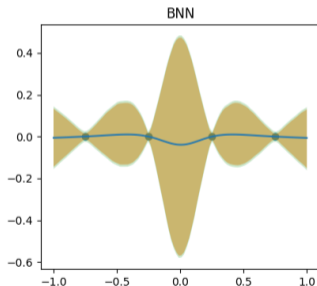
NOT ONLY HARD TO ENCODE FUNCTIONAL PROPERTIES WITH BNNs; SOME PROPERTIES ARE IMPOSSIBLE TO GET

- ▶ For example, a BNN (with RBF activations) is nonstationary in amplitude variance (Williams, 1997)



NOT ONLY HARD TO ENCODE FUNCTIONAL PROPERTIES WITH BNNs; SOME PROPERTIES ARE IMPOSSIBLE TO GET

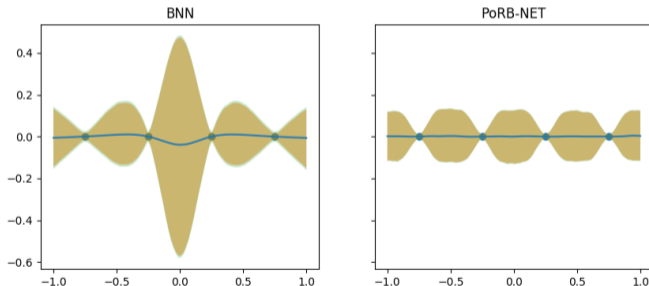
- ▶ For example, a BNN (with RBF activations) is nonstationary in amplitude variance (Williams, 1997)



Question: can we design a Bayesian NN that exhibits stationarity?

NOT ONLY HARD TO ENCODE FUNCTIONAL PROPERTIES WITH BNNs; SOME PROPERTIES ARE IMPOSSIBLE TO GET

- ▶ For example, a BNN (with RBF activations) is nonstationary in amplitude variance (Williams, 1997)



Question: can we design a Bayesian NN that exhibits stationarity? **Yes!**

RELATED WORKS

Expressive priors for Bayesian NNs

- ▶ Functional BNNs (Flam-Shepherd, et.al 2017; Sun et.al, 2019): sample-based optimization w.r.t. reference functional distribution
- ▶ Neural processes (Garnelo et al., 2018): meta-learning to identify functional properties based on many prior examples
- ▶ (Pearce et al., 2019) BNN architectures that recover equivalent GP kernel combinations in the infinite width limit

RELATED WORKS

Expressive priors for Bayesian NNs

- ▶ Functional BNNs (Flam-Shepherd, et.al 2017; Sun et.al, 2019): sample-based optimization w.r.t. reference functional distribution
- ▶ Neural processes (Garnelo et al., 2018): meta-learning to identify functional properties based on many prior examples
- ▶ (Pearce et al., 2019) BNN architectures that recover equivalent GP kernel combinations in the infinite width limit

	user specs	optim. free	finite width	deep
Sun et.al, 2019	yes	no	yes	yes
Garnelo et.al, 2018	no	no	yes	yes
Pearce et.al, 2019	yes	yes	no	yes
PoRB-NET (this work)	yes	yes	yes	not yet

RADIAL BASIS FUNCTION NETWORKS (RBFNs)

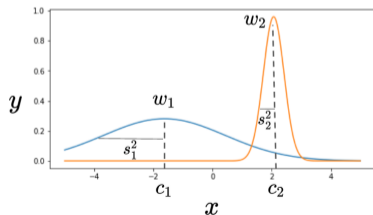
- ▶ Around since the 90s (Gyorfi et.al, 2002), recently renewed attention (Taghi et.al, 2004; Zadeh et.al, 2018)

RADIAL BASIS FUNCTION NETWORKS (RBFNs)

- ▶ Around since the 90s (Gyorfi et.al, 2002), recently renewed attention (Taghi et.al, 2004; Zadeh et.al, 2018)
- ▶ NN based on radial basis $\phi(\cdot)$, e.g., $\phi(x) = \exp(-x^2)$

$$f_{\theta}(x) = b + \sum_{k=1}^K w_k \phi(s_k(x - c_k)),$$

- ▶ $s_k^2 \in \mathbb{R}$: scale
- ▶ $c_k \in \mathbb{R}$: center
- ▶ $w_k \in \mathbb{R}$: output weight
- ▶ $b \in \mathbb{R}$: output bias



COMPARISON RBFN VERSUS BNN FORMULATION (D=1)

$$f_{\theta}(x) = b + \sum_{k=1}^K w_k \phi(s_k(x - c_k))$$

- ▶ $s_k^2 \in \mathbb{R}$: scale
- ▶ $c_k \in \mathbb{R}$: center
- ▶ $w_k \in \mathbb{R}$: output weight
- ▶ $b \in \mathbb{R}$: output bias

$$f_{\theta}(x) = b + \sum_{k=1}^K w_k \phi(v_k x + d_k)$$

- ▶ $v_k \in \mathbb{R}$: input weight
- ▶ $d_k \in \mathbb{R}$: input bias
- ▶ $w_k \in \mathbb{R}$: output weight
- ▶ $b \in \mathbb{R}$: output bias

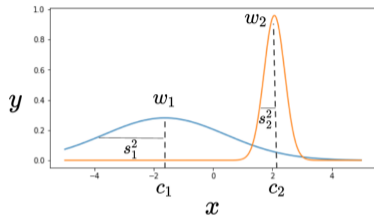
Take-away: priors on different random quantities, RBFN more intuitive

BAYESIAN RBFNS (BARBER ET.AL, 1998)

$$\begin{aligned}
 c_k &\sim \mathcal{N}(0, \sigma_c^2) \\
 s_k^2 &\sim \text{Gamma}(\alpha_s, \beta_s) \\
 w_k &\sim \mathcal{N}(0, \sigma_w^2 I) \\
 B &\sim \mathcal{N}(0, \sigma_b^2) \\
 y_n | x_n, \boldsymbol{\theta} &\sim \mathcal{N}(f_{\boldsymbol{\theta}}(x_n), \sigma_y^2)
 \end{aligned}$$

where

$$f_{\boldsymbol{\theta}}(x) = b + \sum_{k=1}^K w_k \exp\left(-s_k^2(x - c_k)^2\right)$$



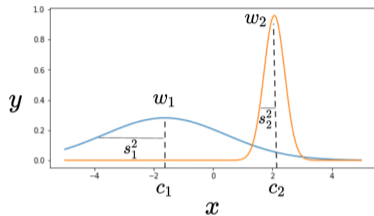
- ▶ $s_k^2 \in \mathbb{R}$: scale
- ▶ $c_k \in \mathbb{R}$: center
- ▶ $w_k \in \mathbb{R}$: output weight
- ▶ $b \in \mathbb{R}$: output bias

BAYESIAN RBFNS (BARBER ET.AL, 1998)

$$\begin{aligned}
 c_k &\sim \mathcal{N}(0, \sigma_c^2) \\
 s_k^2 &\sim \text{Gamma}(\alpha_s, \beta_s) \\
 w_k &\sim \mathcal{N}(0, \sigma_w^2 I) \\
 B &\sim \mathcal{N}(0, \sigma_b^2) \\
 y_n | x_n, \boldsymbol{\theta} &\sim \mathcal{N}(f_{\boldsymbol{\theta}}(x_n), \sigma_y^2)
 \end{aligned}$$

where

$$f_{\boldsymbol{\theta}}(x) = b + \sum_{k=1}^K w_k \exp\left(-s_k^2(x - c_k)^2\right)$$



- ▶ $s_k^2 \in \mathbb{R}$: scale
- ▶ $c_k \in \mathbb{R}$: center
- ▶ $w_k \in \mathbb{R}$: output weight
- ▶ $b \in \mathbb{R}$: output bias

Functional properties still hard or impossible to encode!

FUNCTIONAL PROPERTIES STILL HARD OR IMPOSSIBLE

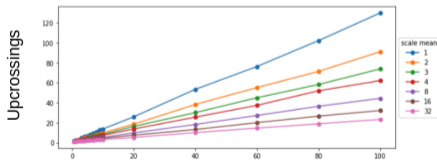
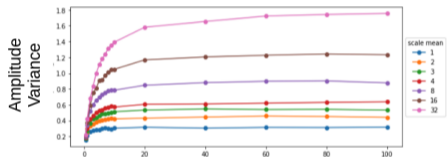
Issues:

- ▶ non-stationary covariance function (Williams, 1997)
- ▶ lengthscale and variance are **coupled**

FUNCTIONAL PROPERTIES STILL HARD OR IMPOSSIBLE

Issues:

- ▶ non-stationary covariance function (Williams, 1997)
- ▶ lengthscale and variance are **coupled**



Density of activations

- ▶ As RBFs concentrate in same region:
 - ▶ summation \implies higher variance
 - ▶ increase in expressivity \implies more upcrossings

POISSON PROCESS RADIAL BASIS FUNCTION NETWORKS (PoRB-NET)

$$c_k \sim \mathcal{N}(0, \sigma_c^2)$$

$$s_k^2 \sim \text{Gamma}(\alpha_s, \beta_s)$$

$$w_k \sim \mathcal{N}(0, \sigma_w^2 I)$$

$$B \sim \mathcal{N}(0, \sigma_b^2)$$

$$y_n | x_n, \boldsymbol{\theta} \sim \mathcal{N}(f_{\boldsymbol{\theta}}(x_n), \sigma_y^2)$$

where

$$f_{\boldsymbol{\theta}}(x) = B + \sum_{k=1}^K w_k \exp(-s_k^2(x - c_k)^2)$$

POISSON PROCESS RADIAL BASIS FUNCTION NETWORKS (PoRB-NET)

$$\mathbf{c} | \lambda \sim \text{Poisson Process}(\lambda)$$

$$s_k^2 \sim \text{Gamma}(\alpha_s, \beta_s)$$

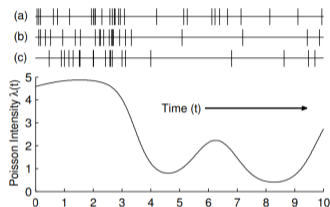
$$w_k \sim \mathcal{N}(0, \tilde{\sigma}_w^2 I)$$

$$B \sim \mathcal{N}(0, \tilde{\sigma}_b^2)$$

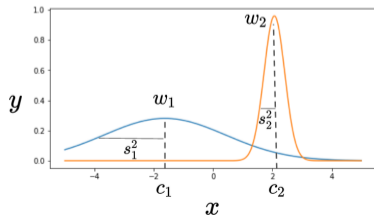
$$y_n | x_n, \boldsymbol{\theta} \sim \mathcal{N}(f_{\boldsymbol{\theta}}(x_n), \sigma_y^2)$$

where

$$f_{\boldsymbol{\theta}}(x) = B + \sum_{k=1}^K w_k \exp\left(-s_k^2(x - c_k)^2\right)$$



(Adams et.al, 2009)



POISSON PROCESS RADIAL BASIS FUNCTION NETWORKS (PoRB-NET)

$$\mathbf{c} \mid \lambda \sim \text{Poisson Process}(\lambda)$$

$$s_k^2 = \lambda^2(c_k)$$

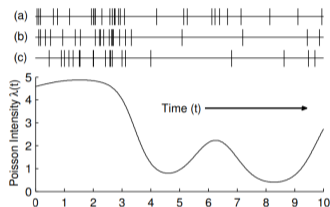
$$w_k \sim \mathcal{N}(0, \tilde{\sigma}_w^2 I)$$

$$B \sim \mathcal{N}(0, \tilde{\sigma}_b^2)$$

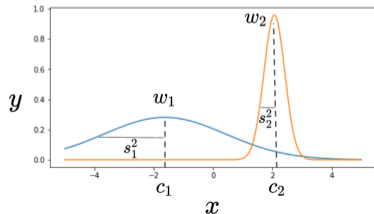
$$y_n \mid x_n, \boldsymbol{\theta} \sim \mathcal{N}(f_{\boldsymbol{\theta}}(x_n), \sigma_y^2)$$

where

$$f_{\boldsymbol{\theta}}(x) = B + \sum_{k=1}^K w_k \exp\left(-s_k^2(x - c_k)^2\right)$$



(Adams et.al, 2009)



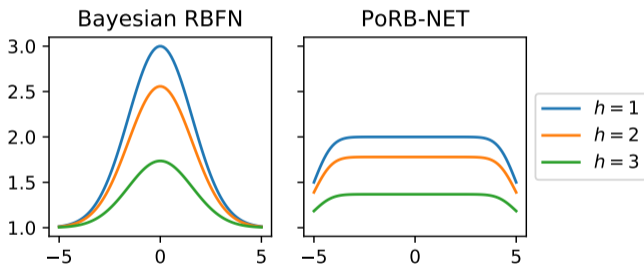
- ▶ We have proposed an expressive prior for NNs

- ▶ We show desirable properties:
 1. Stationarity
 2. Decoupling of lengthscale and amplitude variance
 3. Consistency

- ▶ We demonstrate successful behavior empirically

STATIONARITY

$$\text{Cov}(f(x), f(x+h)) = \sigma_b^2 + \sigma_w^2 \mathbb{E}[K] \underbrace{\mathbb{E}_\theta [\rho(x; \theta) \rho(x+h; \theta)]}_{:=U(x, x+h)}$$



$$U(x_1, x_2) \propto \underbrace{\exp\left(-\frac{(x_1 - x_2)^2}{2(2\sigma_s^2 + \sigma_s^4/\sigma_c^2)}\right)}_{\text{Stationary}} \underbrace{\exp\left(-\frac{x_1^2 + x_2^2}{2(2\sigma_c^2 + \sigma_s^2)}\right)}_{\text{Nonstationary}}$$

$$U(x_1, x_2) = \frac{\lambda}{\Lambda} \sqrt{\frac{\pi}{s^2}} \exp\left\{-s^2 \left(\frac{x_1 - x_2}{2}\right)^2\right\} \left[\Phi((C_1 - x_m)\sqrt{2s^2}) - \Phi((C_0 - x_m)\sqrt{2s^2}\lambda)\right]$$

DECOUPLED LENGTHSCALE AND AMPLITUDE VARIANCE

▶ **Homogeneous Poisson Process**

- ▶ We derive closed-form expression for covariance function
- ▶ Poisson process defined over finite region \mathcal{C}
- ▶ As size of \mathcal{C} tends to infinity,

$$\text{Cov}(f(x_1), f(x_2)) \approx \sigma_b^2 + \tilde{\sigma}_w^2 \exp\left\{-\lambda^2 \left(\frac{x_1 - x_2}{2}\right)^2\right\}$$

▶ **Non-homogeneous Poisson Process**

- ▶ Empirical stationarity

CONSISTENCY

- ▶ Estimator $\hat{g}_n(x)$ is said to be consistent with respect to the true regression function $g_0(x)$ if, as n tends to infinity:

$$\int (\hat{g}_n(x) - g_0(x))^2 dx \xrightarrow{P} 0.$$

- ▶ Posterior consistent over Hellinger neighborhoods if $\forall \epsilon > 0$,

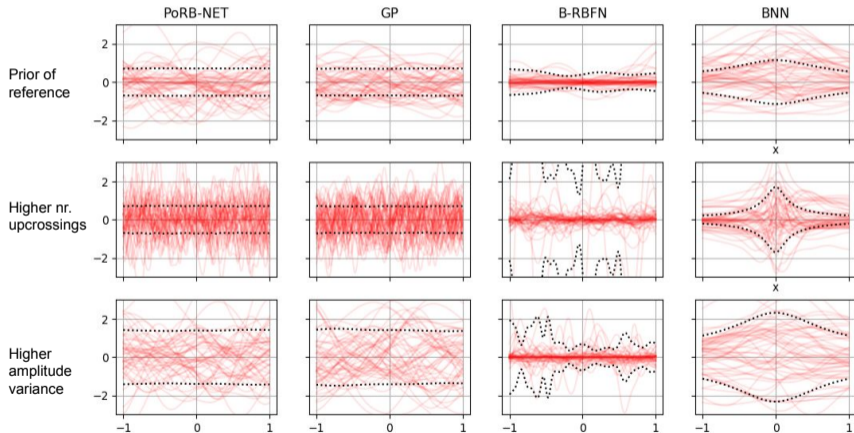
$$p(\{f : D_H(f, f_0) \leq \epsilon\}) \xrightarrow{P} 1.$$

- ▶ (Lee, 2000) shows that Hellinger consistency implies frequentist consistency.

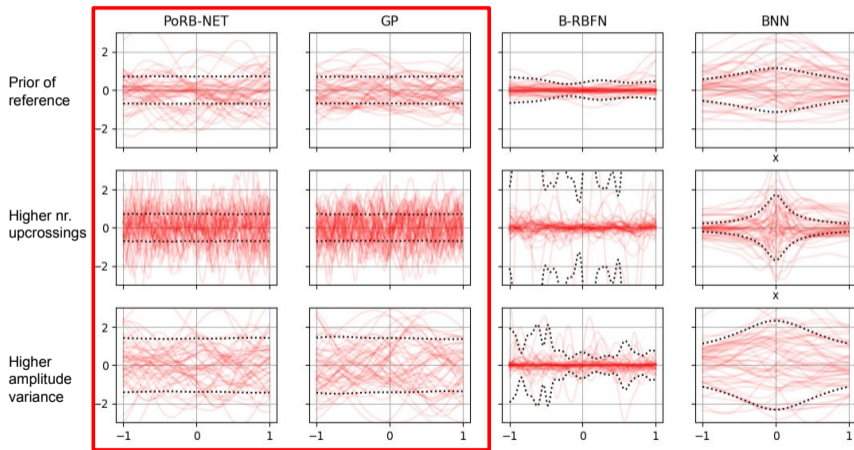
THEOREM (CONSISTENCY OF PoRB-NETs)

A PoRB-NET with uniform intensity function is Hellinger consistent as the number of observations goes to infinity.

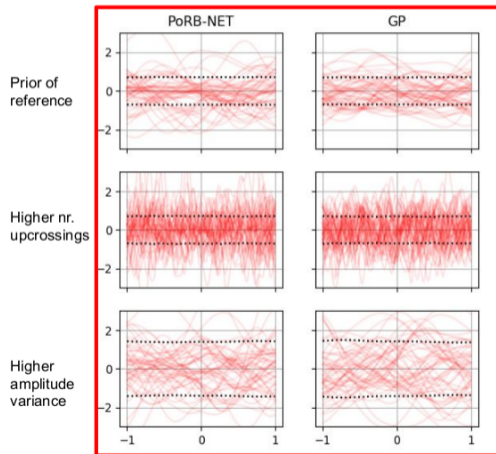
PORB-NET ALLOWS FOR EASY SPECIFICATION OF LENGTHSCALE AND SIGNAL VARIANCE LIKE A GP



PORB-NET ALLOWS FOR EASY SPECIFICATION OF LENGTHSCALE AND SIGNAL VARIANCE LIKE A GP

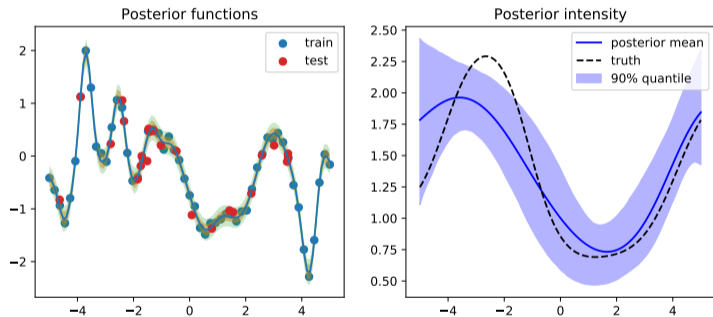


PORB-NET ALLOWS FOR EASY SPECIFICATION OF LENGTHSCALE AND SIGNAL VARIANCE LIKE A GP



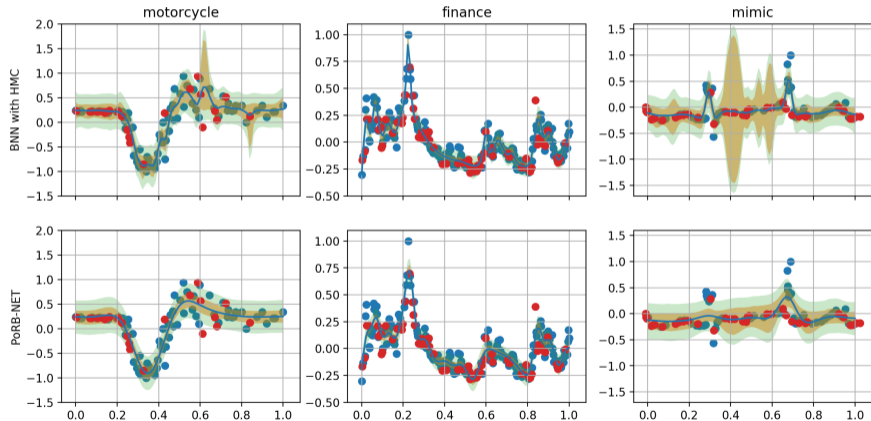
1. stationarity
2. easy specification in a decoupled manner

PORB-NET IS ABLE TO LEARN INPUT-DEPENDENT LENGTHSCALE INFORMATION



PoRB-NET adds more hidden units wherever needed, and adapts architecture width based on the data.

PORB-NET IS ABLE TO CAPTURE NON-STATIONARY PATTERNS IN REAL SCENARIOS, ADAPTING THE LENGTHSCALE LOCALLY



GOAL II: FUNCTIONAL PREDICTION

TAKE-AWAYS...

In this talk, we have...

- ▶ highlighted incapacity of BNNs to express functional properties
- ▶ introduced PoRB-NET, a Bayesian NN prior to encode functional desiderata like a GP
- ▶ proposed an inference scheme to learn input-dependent lengthscale
- ▶ showed theoretical properties: (i) consistency, (ii) decoupling of amplitude and lengthscale
- ▶ validated empirically in synthetic and real datasets

All information online: <https://arxiv.org/abs/1912.05779>

OUTLINE

- ▶ Overview
- ▶ Goal I: Biomarker discovery
- ▶ Goal II: Functional prediction
- ▶ **Wrap-up**

CONCLUSION

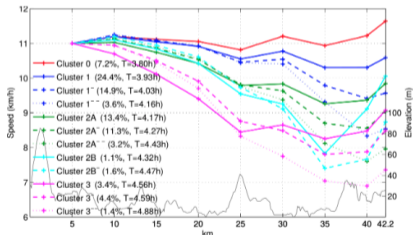
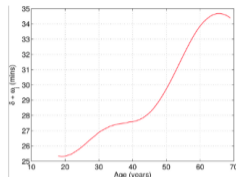
IN THIS TALK...

EXPRESSIVE PRIORS TO ENCODE ASSUMPTIONS/DESIDERATA

1. Structured latent feature model
 - ▶ subpopulation learning
 - ▶ biomarker discovery
2. Novel Bayesian prior for Neural Networks
 - ▶ encoding of stationarity
 - ▶ decoupling of amplitude variance and lengtscale

Other projects...

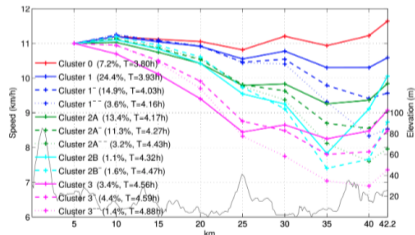
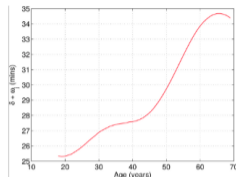
Sport Science



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...

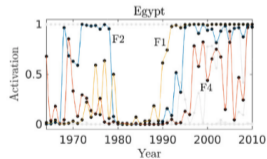
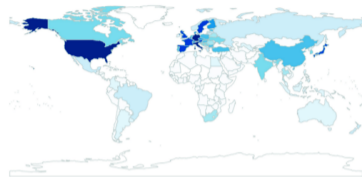
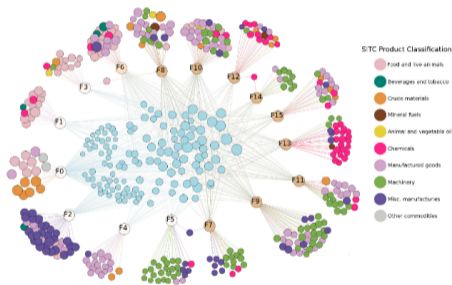
Sport Science



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...

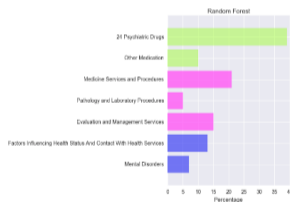
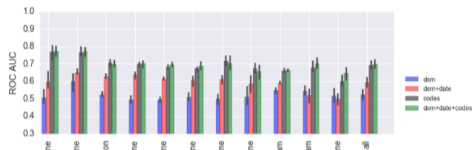
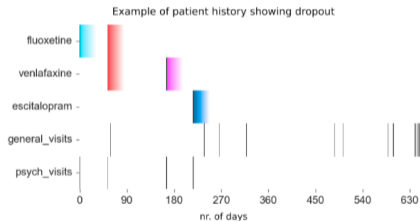
Economics



M. F. Pradier*, Z. Utkovski*, 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.** *Nature Translational Psychiatry*. 2019.

[M. F. Pradier](#), M. Hughes, T. H. McCoy, S. Barroilhet, F. Doshi-Velez and R. H. Perlis. **Predicting Transition from Major Depression to Bipolar Disorder after Antidepressant Initiation.** *Submitted to American Journal of Psychiatry*. 2019.

From the lab to the clinic

- Ongoing user study at MGH, Boston
 - Impact of explanations
 - Usefulness, trust...

Why are these therapies being recommended?

The following **patient features** had the highest contributions to system.13's predictions:



Which antidepressant medication would you be most likely to prescribe in this situation?

Ongoing [M. Jacobs et al 2019]



HARVARD
UNIVERSITY

Patient Details:

Jessica is a 37 year old woman who is married and works part time. She presents with 9 months of depressed mood and lack of appetite. She has a seizure disorder, and current medications include Omeprazole and Celecoxib. Prior treatment with Citalopram had no effect on depressed mood.

System.15 Recommendation: FLUOXETINE

Top 5 therapies with highest probability for stability:



*Stability: continued use of the same medication for at least 3 months

**Dropout: early treatment discontinuation following prescription

Why are these therapies being recommended?

The following **rules** had the highest contributions to system.15's predictions:

1. If *underweight or lack of appetite, favor weight gain, favor Mirtazapine*
2. If *underweight or lack of appetite, avoid appetite suppressants, avoid nausea-inducing, avoid SNRI's, avoid Sertraline*
3. If *lack of response to Paroxetine, avoid SSRI's*

ACKNOWLEDGEMENTS

Special thanks to:

- Beau Coker
- Finale Doshi-Velez
- All members of DTAK!

- Oscar Puig
- Francesca Milletti
- Fernando Perez-Cruz

- Isabel Valera
- Maria Lomeli
- Zoubin Ghahramani



CRCS Center for Research on
Computation and Society

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



HDSI | Harvard Data
Science Initiative



THANK YOU FOR LISTENING!



Looking forward to your questions!
<https://melaniefp.github.io/>

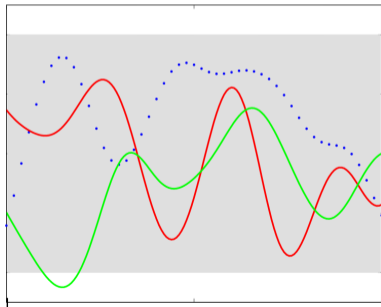
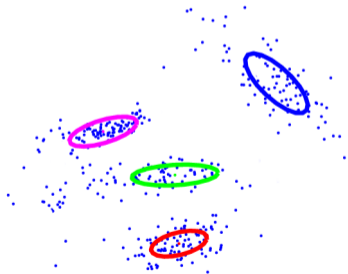
APPENDIX

BAYESIAN NONPARAMETRICS

- ▶ Bayesian: to handle uncertainty

$$p(\text{parameters}|\text{data}) \propto p(\text{data}|\text{parameters})p(\text{parameters})$$

- ▶ Nonparametric: to adapt model complexity (hypothesis generation)

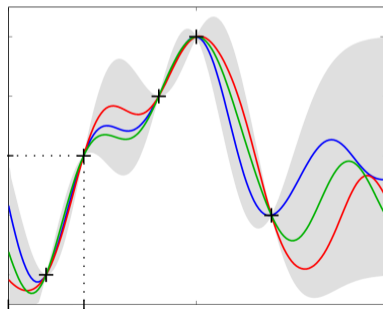
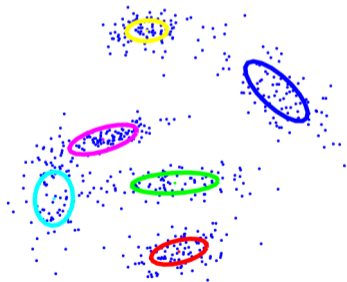


BAYESIAN NONPARAMETRICS

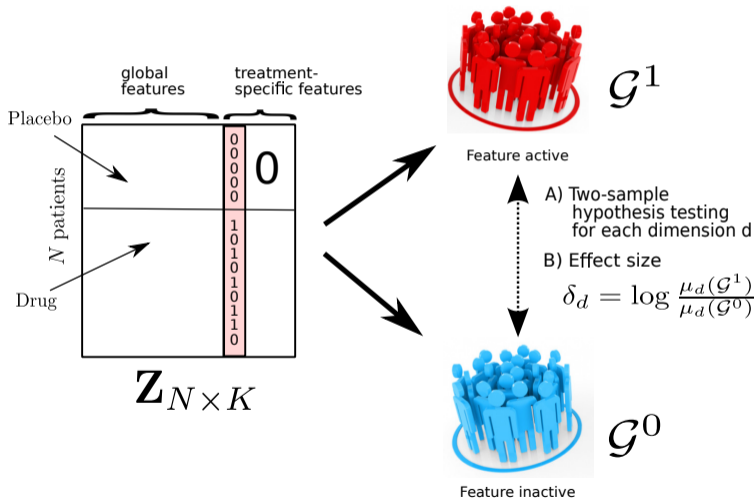
- ▶ Bayesian: to handle uncertainty

$$p(\text{parameters}|\text{data}) \propto p(\text{data}|\text{parameters})p(\text{parameters})$$

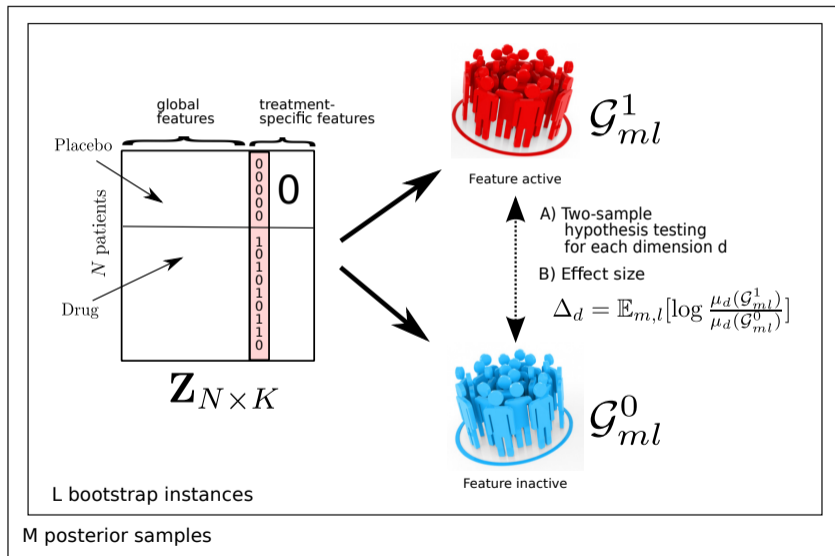
- ▶ Nonparametric: to adapt model complexity (hypothesis generation)



STATISTICAL PROCEDURE FOR BIOMARKER DISCOVERY

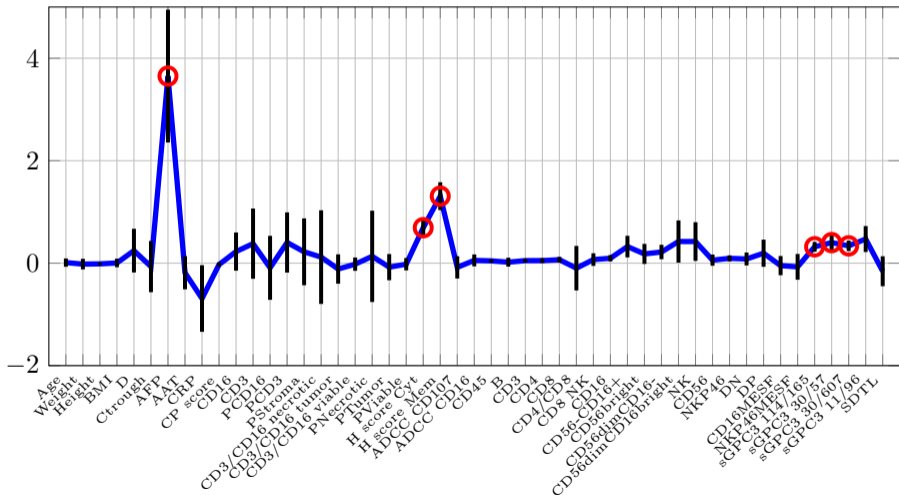


STATISTICAL PROCEDURE FOR BIOMARKER DISCOVERY



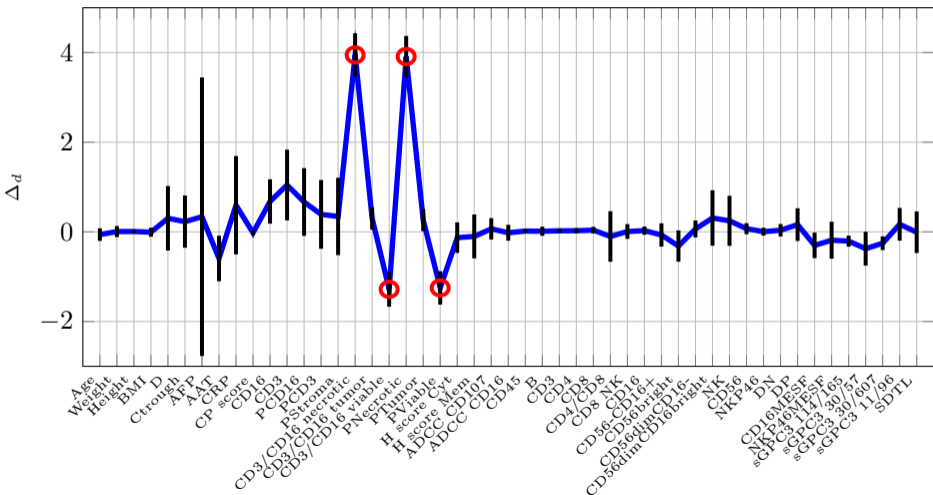
RESULTS: BIOMARKER DISCOVERY

GLOBAL FEATURE F1



RESULTS: BIOMARKER DISCOVERY

GLOBAL FEATURE F2



INDIAN BUFFET PROCESS (IBP)

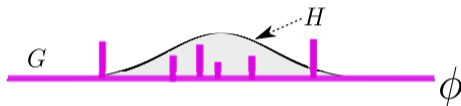
AN ALTERNATIVE CONSTRUCTION

- ▶ underlying block for infinite latent feature models

INDIAN BUFFET PROCESS (IBP)

AN ALTERNATIVE CONSTRUCTION

- ▶ underlying block for infinite latent feature models
- ▶ hierarchy of a Beta process (BP) with multiple Bernoulli processes (BeP)

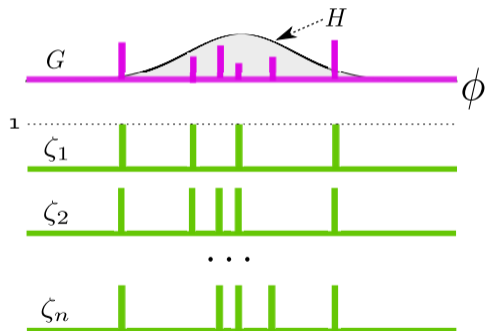


$$G = \sum_{k=1}^{\infty} \pi_k \delta_{\phi_k} \sim \text{BP}(c, \alpha, H)$$

INDIAN BUFFET PROCESS (IBP)

AN ALTERNATIVE CONSTRUCTION

- ▶ underlying block for infinite latent feature models
- ▶ hierarchy of a Beta process (BP) with multiple Bernoulli processes (BeP)



$$G = \sum_{k=1}^{\infty} \pi_k \delta_{\phi_k} \sim \text{BP}(c, \alpha, H)$$

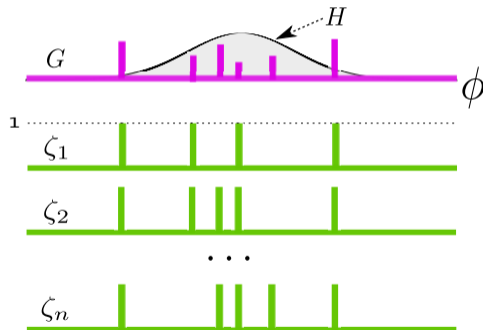
For $n = 1, \dots, \infty$

$$\zeta_n = \sum_{k=1}^{\infty} z_{nk} \delta_{\phi_k} \sim \text{BeP}(G)$$

INDIAN BUFFET PROCESS (IBP)

AN ALTERNATIVE CONSTRUCTION

- ▶ underlying block for infinite latent feature models
- ▶ hierarchy of a Beta process (BP) with multiple Bernoulli processes (BeP)



$$G = \sum_{k=1}^{\infty} \pi_k \delta_{\phi_k} \sim \text{BP}(c, \alpha, H)$$

For $n = 1, \dots, \infty$

$$\zeta_n = \sum_{k=1}^{\infty} z_{nk} \delta_{\phi_k} \sim \text{BeP}(G)$$

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

COMPARISON RBFN VERSUS BNN FORMULATION (D=1)

$$f_{\theta}(x) = B + \sum_{k=1}^K w_k \phi(s_k(x - c_k))$$

$$f_{\theta}(x) = B + \sum_{k=1}^K w_k \phi(v_k x + b_k)$$

$$s_k^2 \sim \text{Gamma}(\alpha_s, \beta_s)$$

$$c_k \sim \mathcal{N}(0, \sigma_c^2)$$

$$w_k \sim \mathcal{N}(0, \sigma_w^2)$$

$$b \sim \mathcal{N}(0, \sigma_0^2)$$

$$v_k^2 \sim \mathcal{N}(0, \sigma_v^2)$$

$$b_k \sim \mathcal{N}(0, \sigma_b^2)$$

$$w_k \sim \mathcal{N}(0, \sigma_w^2)$$

$$b \sim \mathcal{N}(0, \sigma_0^2)$$

Take-away: priors on different random quantities, RBFN more intuitive

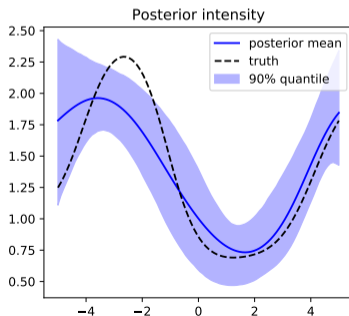
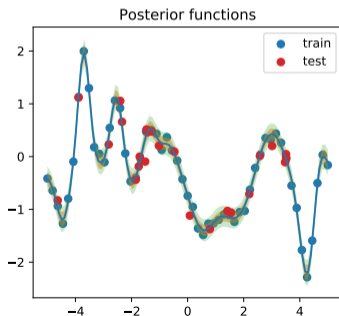
WHAT IF WE DON'T KNOW THE INTENSITY FUNCTION?

Prior on Intensity Function of Poisson Process

$$h \sim \text{GP}(0, C(\cdot, \cdot))$$

$$\lambda^* \sim \text{Gamma}(\alpha_\lambda, \beta_\lambda)$$

$$\lambda(c) = \lambda^* \text{sigmoid}(h(c)),$$



INFERENCE

1. Update network parameters θ given fixed nr. of hidden units K via Hamiltonian Monte Carlo (HMC)

$$p(\theta | \mathbf{y}, \mathbf{x}, K, \lambda) \propto \left(\prod_{n=1}^N \mathcal{N}(y_n; f(x_n; \theta)) \right) \mathcal{N}(b; 0, \sigma_b^2) \left(\prod_{k=1}^K \mathcal{N}(w_k; 0, \sigma_w^2) \lambda(c_k) \right)$$

2. Update network width K via birth/death moves
3. Update point-estimate for Poisson process intensity λ

$$\hat{\lambda}(c) \approx \frac{1}{S} \sum \lambda^* \phi(h^{(s)}(c)),$$

where $h^{(s)} \sim p(h | \mathbf{y}, \mathbf{x}, \theta)$.