

Bayesian Nonparametrics

Part II

Peter Orbanz

1. Constructing nonparametric Bayesian models
 - ▶ Hierarchical and dependent models
 - ▶ Representations
 - ▶ Exchangeability
2. Asymptotics

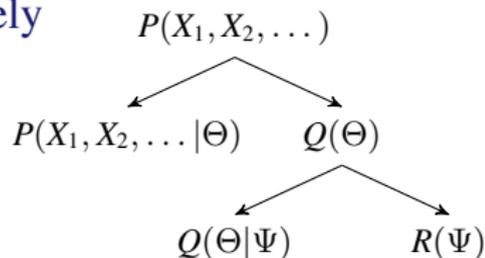
NEW MODELS FROM OLD ONES

HIERARCHICAL MODELS

Apply Bayesian representation recursively

Split parameter Θ :

$$\Theta \rightarrow \Psi \text{ and } \Theta|\Psi$$

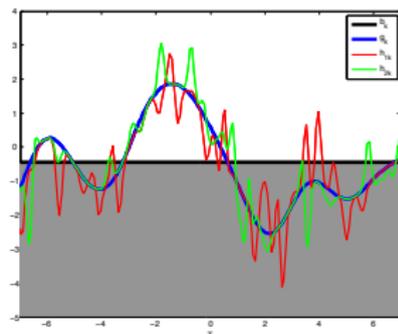


Example: Hierarchical Gaussian process

- ▶ Sample $\Psi \sim R$
(large length-scale, mean 0)
- ▶ Sample $\Theta|\Psi \sim Q(\cdot|\Psi)$
(smaller length scale, mean Ψ)

Decomposes underlying pattern:

- ▶ Low-frequency component Ψ
- ▶ High-frequency component Θ



HIERARCHICAL DIRICHLET PROCESS

Sampling scheme

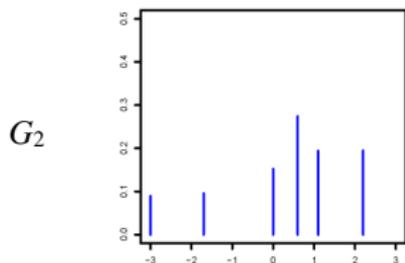
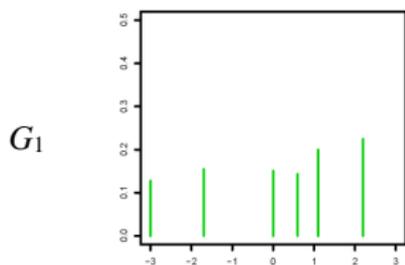
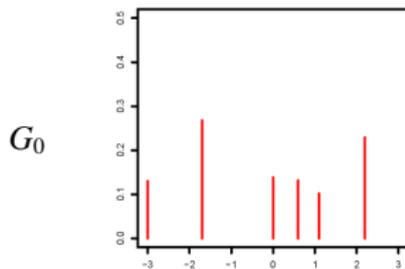
- ▶ Sample $G_0 \sim \text{DP}(\gamma, H)$
- ▶ Sample $G_1, G_2, \dots \sim \text{DP}(\alpha, G_0)$
- ▶ Sample $x_{ij} \sim G_j$

G_1, G_2, \dots have common "vocabulary" of atoms

Application: Nonparametric LDA

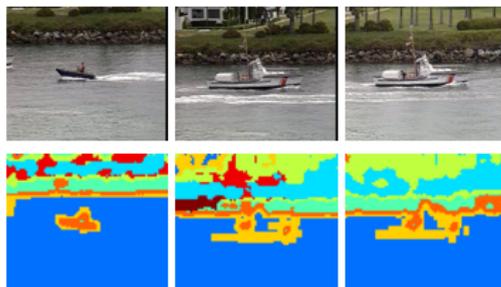
$$G_0 = \sum_{k=1}^{\infty} C_k \delta_{\Theta_k^*} \quad G_j = \sum_{l=1}^{\infty} D_l^j \delta_{\Phi_l^j}$$

- ▶ Θ_k = finite probability (=“topic”)
- ▶ C_k = occurrence probability of topic k
- ▶ Document j drawn from weighted combination of topics, with proportions D_l^j (“admixture model”)



Setting

- ▶ Solution (= pattern) depends on a *covariate*, e.g. time, space,...
- ▶ Example: Video segmentation



For each frame: Solution is a segmentation, i.e. a clustering

Covariate-dependent clustering

$$M(\cdot, t) = \sum_{k=1}^{\infty} C_k(t) \delta_{\Theta_k(t)}(\cdot)$$

for each covariate value t .

DEPENDENT DIRICHLET PROCESS

Dependent Dirichlet process

Model functions $C : T \rightarrow [0, 1]$ and $\Theta : T \rightarrow \Omega_\theta$ with Gaussian processes.

1. Transform GP to have $\text{Beta}(1, \alpha(t))$ marginal distribution for each t .
2. Sample functions $V_1(t), V_2(t), \dots$ from this process.
3. $C_k(t) := V_k(t) \prod_{i=1}^{k-1} (1 - V_i(t))$

Properties

- ▶ Marginal at t is DP $(\alpha(t), G_t)$ with Gaussian base measure G_t .
- ▶ Clustering solutions vary smoothly in t .

Covariate-dependent models: General theme

- ▶ Random object $\Psi \in \Omega_\psi$ with distribution P , covariate space T .
- ▶ Covariate-dependent P : Distribution of random mapping $\hat{\Psi} : T \rightarrow \Omega_\psi$.

EXAMPLES

Applications	Pattern	Bayesian nonparametric model
Classification & regression	Function	Gaussian process
Clustering	Partition	Chinese restaurant process
Density estimation	Density	Dirichlet process mixture
Hierarchical clustering	Hierarchical partition	Dirichlet/Pitman-Yor diffusion tree, Kingman's coalescent, Nested CRP
Latent variable modelling	Features	Beta process/Indian buffet process
Survival analysis	Hazard	Beta process, Neutral-to-the-right process
Power-law behaviour		Pitman-Yor process, Stable-beta process
Dictionary learning	Dictionary	Beta process/Indian buffet process
Dimensionality reduction	Manifold	Gaussian process latent variable model
Deep learning	Features	Cascading/nested Indian buffet process
Topic models	Atomic distribution	Hierarchical Dirichlet process
Time series		Infinite HMM
Sequence prediction	Conditional probs	Sequence memoizer
Reinforcement learning	Conditional probs	infinite POMDP
Spatial modelling	Functions	Gaussian process, dependent Dirichlet process
Relational modelling		Infinite relational model, infinite hidden relational model, Mondrian process
...

REPRESENTATIONS

Densities

$$P(dx) = p(x)\lambda(dx) \quad P(A) = \int_A p(x)\lambda(dx)$$

We call λ the *carrier measure* and p the *density* of P w.r.t. λ .

Useful carrier measures

- ▶ λ should be translation-invariant.
- ▶ Such measures exist only on certain spaces, roughly speaking:
On finite-dimensional spaces.

Consequence: Representation problem 1

- ▶ Nonparametric models: No useful carrier measure on parameter space.
- ▶ We have to find alternatives to density representation.

THE BAYES EQUATION

Bayesian model: General case

Prior distribution Q , likelihood $P[X \in \cdot | \Theta]$, posterior $Q[\Theta \in \cdot | X = x]$

Bayes' Theorem

If the posterior has a density w.r.t. the prior for each x , then

$$Q[d\theta | X = x] = \frac{dQ[\cdot | X = x]}{dQ(\cdot)} Q(d\theta) = \frac{dP[X \in \cdot | \theta]}{dP(X \in \cdot)}(x) Q(d\theta)$$

The “Bayes equation” is a density of the posterior with respect to the prior.

Representation Problem 2

- ▶ For many nonparametric models, this density cannot exist for all x .
- ▶ Such models are called *undominated*.
- ▶ Random discrete measure models are generally undominated.

In other words:

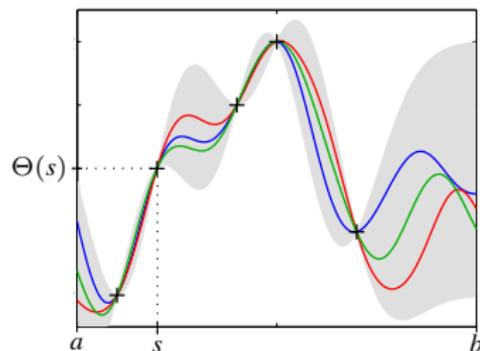
NPB models do not generally satisfy Bayes' theorem.

GAUSSIAN PROCESSES

Nonparametric regression

Patterns = continuous functions, say on $[a, b]$:

$$\theta : [a, b] \rightarrow \mathbb{R} \quad \mathcal{T} = C[a, b]$$



Recall definition

$$\Theta \sim \text{GP} \quad \Leftrightarrow \quad (\Theta(s_1), \dots, \Theta(s_d)) \quad \text{is } d\text{-dimensional Gaussian}$$

for any finite set $S \subset [a, b]$.

Construction: Intuition

- ▶ The marginal of the GP for any finite $S \subset [a, b]$ is a Gaussian.
- ▶ All these Gaussians are marginals of each other.
- ▶ Conversely: If we start with such Gaussians for all S , do they define a GP?

They do. The theorems which guarantee this are called *extension theorems* or *projective limit theorems*.

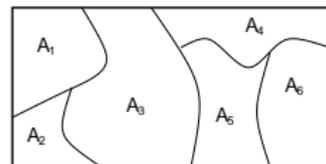
CONSTRUCTING RANDOM MEASURES

Idea

- ▶ GP: We have constructed a *random function* Θ .
- ▶ If Θ is a *random measure*, can we construct it in a similar way?

Extension theorem

- ▶ For a finite partition $I = (A_1, \dots, A_d)$ of V , suppose we know the distribution P_I of $(\Theta(A_1), \dots, \Theta(A_d))$.
- ▶ If the P_I for all partitions I are projective (= are marginals of each other), they define a unique random measure Θ on V .



Example: DP

Choose P_I as Dirichlet distribution with parameters α and $(G_0(A_1), \dots, G_0(A_d))$.
Then $\Theta \sim \text{DP}(\alpha, G_0)$.

Stick-breaking

- ▶ Simple; most widely used where applicable.
- ▶ Constructive.
- ▶ Available only for few models (DP, Pitman-Yor process, normalized inverse Gaussian process, beta process).

Projective limits

- ▶ Generally applicable.
- ▶ Mathematically more challenging, many open problems.

Representations by known stochastic processes

- ▶ E.g. Lévy process and Poisson process representations.
- ▶ Often come with a useful set of theoretical tools.

Conjugate models

- ▶ How can we compute a posterior without a Bayes equation?
- ▶ Virtually all NPB models (DP, GP, etc) are conjugate.

Functional vs structural conjugacy

Functional conjugacy: There is a mapping

prior hyperparameter \times data \mapsto posterior hyperparameter

Structural conjugacy: Closure under sampling, but no functional conjugacy.

Example

Neutral-to-the-right processes are structurally but not functionally conjugate.

Constructing conjugate models

- ▶ In hierarchical models: Use conjugate components.
- ▶ Roughly: Projective limits of fct. conjugate marginals are fct. conjugate.

EXCHANGEABILITY

MOTIVATION

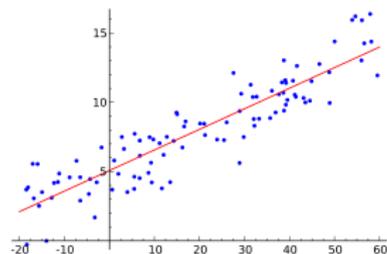
Can we justify our assumptions?

Recall:

$$\text{data} = \text{pattern} + \text{noise}$$

In Bayes' theorem:

$$Q(d\theta|x_1, \dots, x_n) = \frac{\prod_{j=1}^n p(x_j|\theta)}{p(x_1, \dots, x_n)} Q(d\theta)$$



Exchangeability

X_1, X_2, \dots are *exchangeable* if $P(X_1, X_2, \dots)$ is invariant under any permutation σ :

$$P(X_1 = x_1, X_2 = x_2, \dots) = P(X_1 = x_{\sigma(1)}, X_2 = x_{\sigma(2)}, \dots)$$

In words:

Order of observations does not matter.

De Finetti's Theorem

$$P(X_1 = x_1, X_2 = x_2, \dots) = \int_{M(\mathcal{X})} \left(\prod_{j=1}^{\infty} \theta(X_j = x_j) \right) Q(d\theta)$$

\Updownarrow

X_1, X_2, \dots exchangeable

where:

- ▶ $M(\mathcal{X})$ is the set of probability measures on \mathcal{X}
- ▶ θ are values of a random probability measure Θ with distribution Q

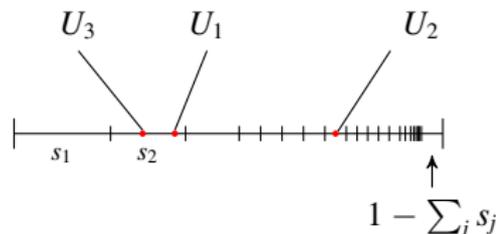
Implications

- ▶ Exchangeable data decomposes into pattern and noise
- ▶ More general than i.i.d.-assumption
- ▶ Caution: θ is in general an ∞ -dimensional quantity

EXCHANGEABILITY: RANDOM PARTITIONS

Paint-box distribution

- ▶ Weights $s_1, s_2, \dots \geq 0$ with $\sum s_j \leq 1$
- ▶ $U_1, U_2, \dots \sim \text{Uniform}[0, 1]$



Random partition of \mathbb{N} :

$i, j \in \mathbb{N}$ in same block $\Leftrightarrow U_i, U_j$ in same interval

$\{i\}$ separate block $\Leftrightarrow U_i$ in interval $1 - \sum s_j$

Kingman's Theorem

Random partition π of \mathbb{N} exchangeable

\Updownarrow

Mixture of paint-boxes $\beta(\cdot | \mathbf{s})$: $P(\pi) = \int \beta(\pi | \mathbf{s}) Q(d\mathbf{s})$

EXCHANGEABILITY: RANDOM GRAPHS

Random graph with independent edges

Given: $\theta : [0, 1]^2 \rightarrow [0, 1]$ symmetric function

- ▶ $U_1, U_2, \dots \sim \text{Uniform}[0, 1]$
- ▶ Edge (i, j) present:

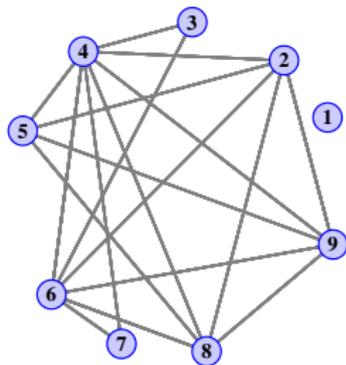
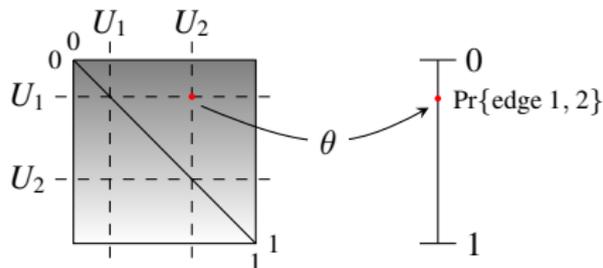
$$(i, j) \sim \text{Bernoulli}(\theta(U_i, U_j))$$

Call this distribution $P(\mathcal{G}|\theta)$.

Aldous-Hoover Theorem

Random graph \mathcal{G} exchangeable

$$\begin{aligned} &\Updownarrow \\ P(\mathcal{G}) &= \int_{\mathcal{T}} P(\mathcal{G}|\theta) Q(d\theta) \end{aligned}$$



GENERAL THEME: SYMMETRY

Other types of exchangeable data

Data	Theorem	Mixture of...	Applications
Points	de Finetti	I.i.d. point sequences	“Standard” models
Sequences	Diaconis-Freedman	Markov chains	Time series
Partition	Kingman	"Paint-box" partitions	Clustering
Graphs	Aldous-Hoover	Graphs with independent edges	Networks
Arrays	Aldous-Hoover	Arrays with independent entries	Collaborative filtering

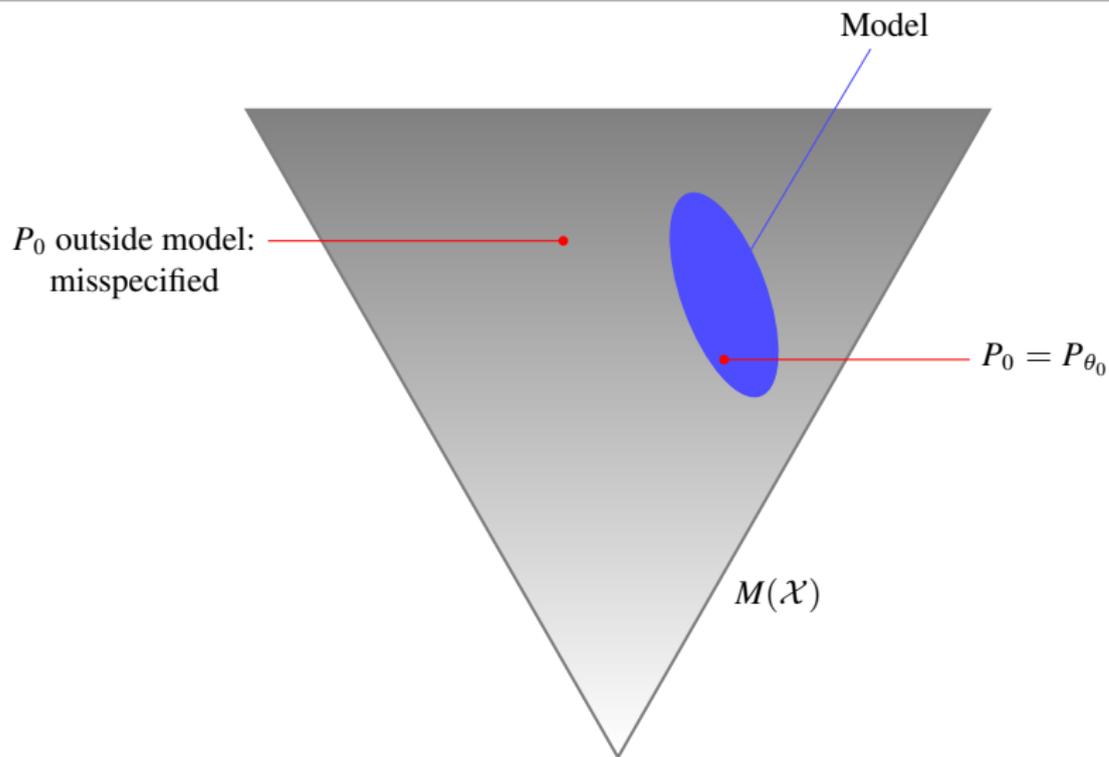
Ergodic decomposition theorems

$$\mu(X) = \int_{\Omega} \mu[X|\Phi = \phi] \nu(\phi)$$

- ▶ Symmetry (group invariance) on lhs \longrightarrow Integral decomposition on rhs
- ▶ Permutation invariance on lhs \longrightarrow Independence on rhs

ASYMPTOTICS

SUPPORT OF PRIORS



SUPPORT OF NONPARAMETRIC PRIORS

Large support

- ▶ Support of nonparametric priors is larger (∞ -dimensional) than of parametric priors (finite-dimensional).
- ▶ However: No uniform prior (or even “neutral” improper prior) exists on $M(\mathcal{X})$.

Interpretation of nonparametric prior assumptions

Concentration of nonparametric prior on subset of $M(\mathcal{X})$ typically represents structural prior assumption.

- ▶ GP regression with unknown bandwidth:
 - ▶ Any continuous function possible
 - ▶ Prior can express e.g. “very smooth functions are more probable”
- ▶ Clustering: Expected number of clusters is...
 - ▶ ...small \longrightarrow CRP prior
 - ▶ ...power law \longrightarrow two-parameter CRP

POSTERIOR CONSISTENCY

Definition 1 (weak consistency of Bayesian models)

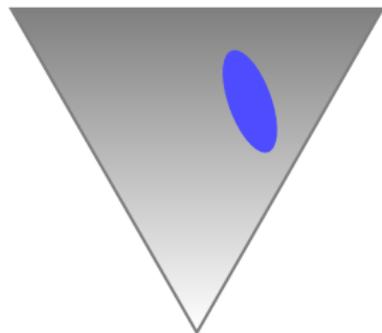
Suppose we sample $P_0 = P_{\theta_0}$ from the prior and generate data from P_0 . If the posterior converges to δ_{θ_0} for $n \rightarrow \infty$ *with probability one under the prior*, the model is called *consistent*.

Doob's Theorem

Under very mild conditions, Bayesian models are consistent in the weak sense.

Problem

- ▶ Definition holds up to a set of probability zero under the prior.
- ▶ This set can be huge and is a prior assumption.



Definition 2 (frequentist consistency of Bayesian models)

A Bayesian model is *consistent at P_0* if the posterior converges to δ_{P_0} with growing sample size.

CONVERGENCE RATES

Objective

How quickly does posterior concentrate at θ_0 as $n \rightarrow \infty$?

Measure: Convergence rate

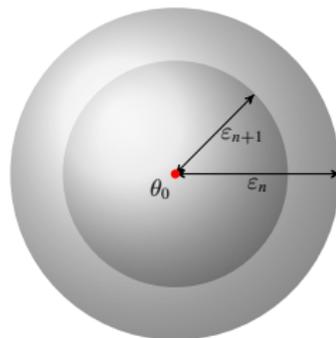
- ▶ Find smallest balls $B_{\varepsilon_n}(\theta_0)$ for which

$$Q(B_{\varepsilon_n}(\theta_0) | X_1, \dots, X_n) \xrightarrow{n \rightarrow \infty} 1$$

- ▶ Rate = sequence $\varepsilon_1, \varepsilon_2, \dots$

The best we can hope for

- ▶ Optimal rate is $\varepsilon_n \propto n^{-1/2}$
- ▶ Given by optimal convergence of estimators
- ▶ Achieved in smooth parametric models



Technical tools

Sieves, covering number, metric entropies... \longrightarrow familiar from learning theory!

Consistency

- ▶ DP mixtures: Consistent in many cases. No blanket statements.
- ▶ Range of consistency results for GP regression

Convergence rates: Example

Bandwidth adaptation with GPs:

- ▶ True parameter $\theta_0 \in C^\alpha[0, 1]^d$, smoothness α unknown
- ▶ With gamma prior on GP bandwidth:

Convergence rate is $n^{-\alpha/(2\alpha+d)}$

Bernstein-von Mises Theorems

- ▶ Class of theorems establishing that posterior is asymptotically normal.
- ▶ Available for Gaussian processes and various regression settings.

REFERENCES I

- [Ald81] David J. Aldous. Representations for Partially Exchangeable Arrays of Random Variables. *Journal of Multivariate Analysis*, 11:581–598, 1981.
- [Fer73] T. S. Ferguson. A Bayesian analysis of some nonparametric problems. *Ann. Statist.*, 1(2), 1973.
- [Gho10] S. Ghosal. Dirichlet process, related priors and posterior asymptotics. In N. L. Hjort et al., editors, *Bayesian Nonparametrics*, pages 36–83. Cambridge University Press, 2010.
- [GvdV07] Subhashis Ghosal and Aad van der Vaart. Posterior convergence rates of Dirichlet mixtures at smooth densities. *Ann. Statist.*, 35(2):697–723, 2007.
- [Kal05] Olav Kallenberg. *Probabilistic Symmetries and Invariance Principles*. Springer, 2005.
- [Kin78] J. F. C. Kingman. The representation of partition structures. *Journal of the London Mathematical Society*, 18:374–380, 1978.
- [KvdV06] B. J. K. Kleijn and A. W. van der Vaart. Misspecification in infinite-dimensional Bayesian statistics. *Annals of Statistics*, 34(2):837–877, 2006.
- [LP10] A. Lijoi and I. Prünster. Models beyond the Dirichlet process. In N. L. Hjort, C. Holmes, P. Müller, and S. G. Walker, editors, *Bayesian Nonparametrics*. Cambridge University Press, 2010.
- [Mac00] S. N. MacEachern. Dependent Dirichlet processes. Technical report, Ohio State University, 2000.
- [Orb09] P. Orbanz. Construction of nonparametric Bayesian models from parametric bayes equations. In *Advances in Neural Information Processing Systems (NIPS)*, 2009.
- [Orb11] P. Orbanz. Projective limit random probabilities on Polish spaces. *Electronic Journal of Statistics*, 5:1354–1373, 2011.
- [Orb12] P. Orbanz. Nonparametric priors on complete separable metric spaces. 2012.
- [RT09] D.M. Roy and Y.-W. Teh. The Mondrian process, 2009.
- [Sch65] L. Schwartz. On Bayes procedures. *Z. Wahr. Verw. Gebiete*, 4:10–26, 1965.
- [Sch95] M. J. Schervish. *Theory of Statistics*. Springer, 1995.
- [TJ10] Y. W. Teh and M. I. Jordan. Hierarchical Bayesian nonparametric models with applications. In N. L. Hjort, C. Holmes, P. Müller, and S. G. Walker, editors, *Bayesian Nonparametrics*. Cambridge University Press, 2010.

REFERENCES II

- [TJBB06] Y. W. Teh, M. I. Jordan, M. J. Beal, and D. M. Blei. Hierarchical Dirichlet processes. *J. Amer. Statist. Assoc.*, (476):1566–1581, 2006.
- [vdV98] A. van der Vaart. *Asymptotic Statistics*. Cambridge University Press, 1998.
- [vdVvZ08a] A. W. van der Vaart and J. H. van Zanten. Rates of contraction of posterior distributions based on Gaussian process priors. *Ann. Statist.*, 36(3):1435–1463, 2008.
- [vdVvZ08b] A. W. van der Vaart and J. H. van Zanten. Reproducing kernel Hilbert spaces of Gaussian priors. In *Pushing the limits of contemporary statistics: contributions in honor of Jayanta K. Ghosh*, volume 3 of *Inst. Math. Stat. Collect.*, pages 200–222. Inst. Math. Statist., Beachwood, OH, 2008.