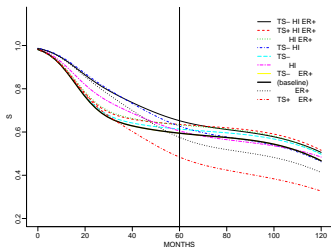


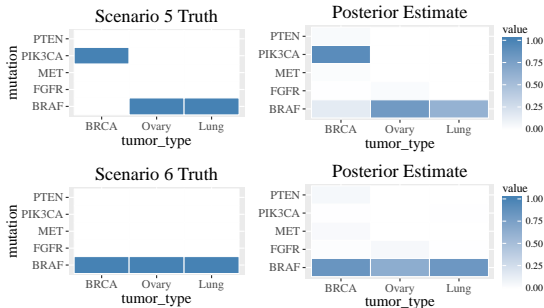
Nonparametric Bayesian Data Analysis for Causal Inference

Part 2 – Regression

PETER MÜLLER, UT Austin



regression



left = truth; right = estimate as $p(a)$ over repeat sim.

Slides: www.math.utexas.edu/users/pmueller/osu.pdf

2. Regression

Regression: $y_i \mid x_i = x \sim F_x(y_i)$.

2. Regression

Regression: $y_i | x_i = x \sim F_x(y_i)$.

1. NP on residual: $y_i = f_\theta(x_i) + \epsilon_i$, $\epsilon_i \sim G$ and $G \sim p(G)$.

Semiparametric Bayes, density estimation for residuals ϵ_i , e.g., PT prior (Hanson & Johnson, 2002 JASA).

2. Regression

Regression: $y_i | x_i = x \sim F_x(y_i)$.

1. NP on residual: $y_i = f_\theta(x_i) + \epsilon_i$, $\epsilon_i \sim G$ and $G \sim p(G)$.

Semiparametric Bayes, density estimation for residuals ϵ_i , e.g., PT prior (Hanson & Johnson, 2002 JASA).

2. Random regression mean function :

$$y_i = f(x_i) + \epsilon_i \text{ and } f(\cdot) \sim p(f)$$

GP prior, wavelet bases, neural networks, hierarchical mixture of experts, etc.

2. Regression

Regression: $y_i | x_i = x \sim F_x(y_i)$.

1. NP on residual: $y_i = f_\theta(x_i) + \epsilon_i$, $\epsilon_i \sim G$ and $G \sim p(G)$.

Semiparametric Bayes, density estimation for residuals ϵ_i , e.g., PT prior (Hanson & Johnson, 2002 JASA).

2. Random regression mean function :

$$y_i = f(x_i) + \epsilon_i \text{ and } f(\cdot) \sim p(f)$$

GP prior, wavelet bases, neural networks, hierarchical mixture of experts, etc.

3. Fully non-parametric regression:

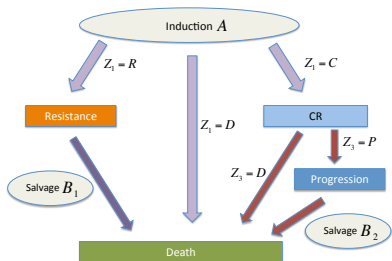
$$y_i | x_i \sim F_{x_i}, \text{ with } \mathcal{F} = \{F_x, x \in X\} \sim p(\mathcal{F}).$$

For example, DDP model, dependent PT etc.

Introduce the DDP next ...

Example 1: Dynamic treatment regimen

Xu et al. (2016 JASA)



Problem: Frontline therapy (A) is randomized, salvage therapy (B) is usually **not randomized**. Adjust for the lack of randomization.

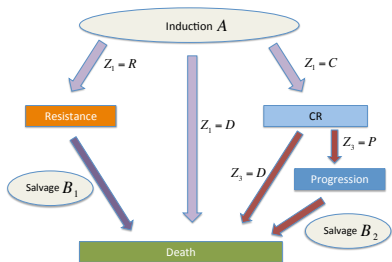
Motivating leukemia trial

Aim: BNP approach to evaluate DTRs, using model-based inference to undo the lack of randomization.

- 4 **induction** trts: FAI, FAI+ATRA, FAI+GCSF, FAI+ATRA+GCSF.
- 2 **salvage** trts: HDAC or not.

Example 1: Dynamic treatment regimen

Xu et al. (2016 JASA)



Problem: Frontline therapy (A) is randomized, salvage therapy (B) is usually **not randomized**. Adjust for the lack of randomization.

Motivating leukemia trial

Aim: BNP approach to evaluate DTRs, using model-based inference to undo the lack of randomization.

- 4 **induction** trts: FAI, FAI+ATRA, FAI+GCSF, FAI+ATRA+GCSF.
- 2 **salvage** trts: HDAC or not.

BNP Model for Evaluating DTRs

Data:

Outcome: $Y^k = \log(T^k) = (\log) k^{th}$ transition time (e.g., R \rightarrow D)

Covariates: \mathbf{x}^k , incl. T^ℓ , $\ell < k$

BNP Model for Evaluating DTRs

Data:

Outcome: $Y^k = \log(T^k) = (\log) k^{\text{th}}$ transition time (e.g., R \rightarrow D)

Covariates: \mathbf{x}^k , incl. T^ℓ , $\ell < k$

Pars: $\mathcal{F} = \{F^k; k = 1, \dots, K\}$, (unknown) distributions of 7 transition times

Likelihood:

$$\prod_{k=1}^K p(Y^k | \mathbf{x}^k, \mathcal{F}) = \prod_{k=1}^K F_{\mathbf{x}^k}^k(Y^k)$$

BNP Model for Evaluating DTRs

Data:

Outcome: $Y^k = \log(T^k) = (\log) k^{\text{th}}$ transition time (e.g., $R \rightarrow D$)

Covariates: \mathbf{x}^k , incl. T^ℓ , $\ell < k$

Pars: $\mathcal{F} = \{F^k; k = 1, \dots, K\}$, (unknown) distributions of 7 transition times

Likelihood:

$$\prod_{k=1}^K p(Y^k | \mathbf{x}^k, \mathcal{F}) = \prod_{k=1}^K F_{\mathbf{x}^k}^k(Y^k)$$

Prior: BNP prior for \mathcal{F}

$$\mathcal{F} = \{F_{\mathbf{x}}^k; \mathbf{x} \in X, \} \sim \text{DDP}, \quad k = 1, \dots, K$$

with $F_{\mathbf{x}}^k = \sum_{h=0}^{\infty} p_h^k N(y; \theta_{h,\mathbf{x}}^k, \sigma^k)$.

BNP Model for Evaluating DTRs

Data:

Outcome: $Y^k = \log(T^k) = (\log) k^{\text{th}}$ transition time (e.g., $R \rightarrow D$)

Covariates: \mathbf{x}^k , incl. T^ℓ , $\ell < k$

Pars: $\mathcal{F} = \{F^k; k = 1, \dots, K\}$, (unknown) distributions of 7 transition times

Likelihood:

$$\prod_{k=1}^K p(Y^k | \mathbf{x}^k, \mathcal{F}) = \prod_{k=1}^K F_{\mathbf{x}^k}^k(Y^k)$$

Prior: BNP prior for \mathcal{F}

$$\mathcal{F} = \{F_{\mathbf{x}}^k; \mathbf{x} \in X, \} \sim \text{DDP}, \quad k = 1, \dots, K$$

with $F_{\mathbf{x}}^k = \sum_{h=0}^{\infty} p_h^k N(y; \theta_{h,\mathbf{x}}^k, \sigma^k)$.

GP prior on $\{\theta_{h,\mathbf{x}}^k\}_{\mathbf{x}}$

Prior: (skip “ k ” superindex for the moment)

$$F_x = \sum_{h=0}^{\infty} p_h N(y; \theta_{h,x}^k, \sigma).$$

Prior: (skip “ k ” superindex for the moment)

$$F_x = \sum_{h=0}^{\infty} p_h N(y; \theta_{h,x}^k, \sigma).$$

- stick-breaking prior on p_h
- GP prior on the functions $\{\theta_{h,x}\}_x$, dependent **across** x , independent across h

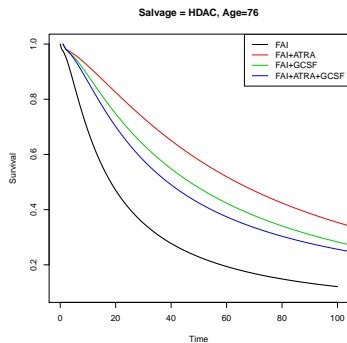
Results: Survival regression and optimal policy

Survival regression: for each T_x^k , using
DDP mix of normal

Results: Survival regression and optimal policy

Survival regression: for each T_x^k , using
DDP mix of normal

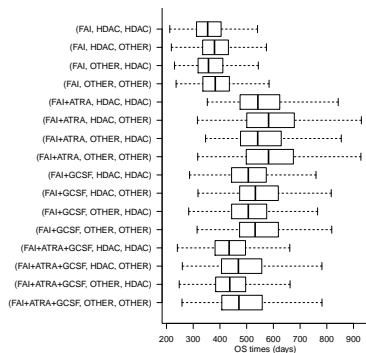
Prior support: full prior support;
BNP is *always right*; this mitigates
concerns about extrapolation.



survival regr for T^{PD}

Comparing policies

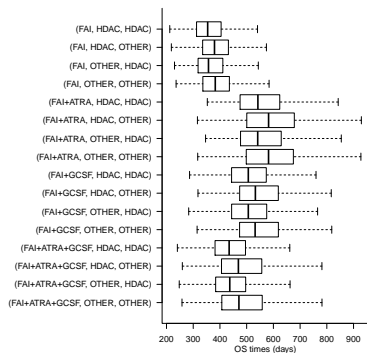
Overall survival for alternative policies (A , B_1 , B_2).



Potential outcomes: evaluate mean OS for possible treatment policies

Comparing policies

Overall survival for alternative policies (A , B_1 , B_2).

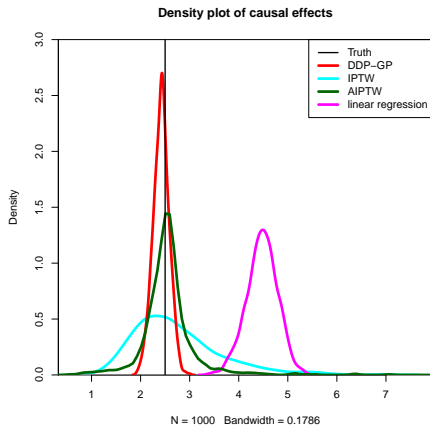


Potential outcomes: evaluate mean OS for possible treatment policies

Optimal policy: compare by mean OS

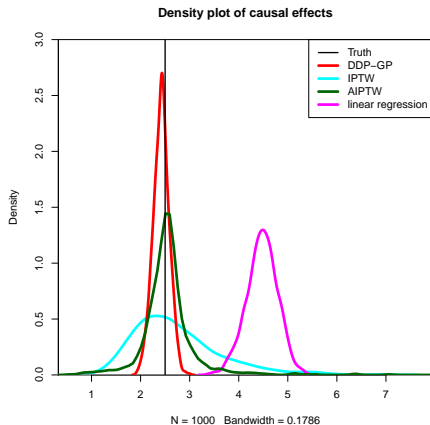
Comparison with double robust methods

Two simulations to compare with inverse prob weighting, using correct model (left) and mis-specified model (right)

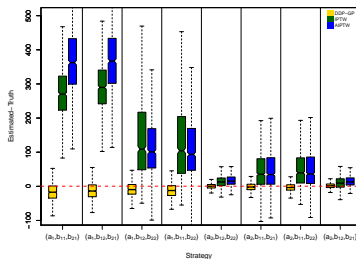


Comparison with double robust methods

Two simulations to compare with inverse prob weighting, using correct model (left) and mis-specified model (right)



single event time
(correct model)



DTR, with both
models wrong

Example 2: Semicompeting risks

Xu, Scharfstein, M and Daniels (2019, arXiv). Another application of (almost) the same model for pairs of event times.

Event times: progression P_j & overall survival D_j
under control ($j = 0$) and treatment ($j = 1$).

Example 2: Semicompeting risks

Xu, Scharfstein, M and Daniels (2019, arXiv). Another application of (almost) the same model for pairs of event times.

Event times: progression P_j & overall survival D_j
under control ($j = 0$) and treatment ($j = 1$).

Censoring: D_j censors P_j ;
and independent censoring C_j

Example 2: Semicompeting risks

Xu, Scharfstein, M and Daniels (2019, arXiv). Another application of (almost) the same model for pairs of event times.

Event times: progression P_j & overall survival D_j
under control ($j = 0$) and treatment ($j = 1$).

Censoring: D_j censors P_j ;
and independent censoring C_j

Inference: compare P_j adjusting for D_j

Inference target: conditional odds

$$\tau_x(u) = \frac{p_x(P_1 < u \mid D_0 > u, D_1 > u)}{p_x(P_0 < u \mid D_0 > u, D_1 > u)}$$

Joint distribution for P_j, D_j

Identifiability: Let

$$G_j = p(D_j)$$

Joint distribution for P_j, D_j

Identifiability: Let

$$G_j = p(D_j)$$

and

$$V_j(s | t) = p(P_j \leq s, P_j < D_j | D_j = t).$$

$s < t$ (for the moment, ignoring regression on “ \mathbf{x} ”).

Under random censoring G_j and V_j are identifiable –
just use the corresponding sample statistics.

Joint distribution for P_j, D_j

Identifiability: Let

$$G_j = p(D_j)$$

and

$$V_j(s | t) = p(P_j \leq s, P_j < D_j | D_j = t).$$

$s < t$ (for the moment, ignoring regression on “ \mathbf{x} ”).

Under random censoring G_j and V_j are identifiable – just use the corresponding sample statistics.

Bivariate sub-distribution: together G_j & V_j define

$$\tilde{F}_1(s, t) = p(P_1 \leq s, D_1 \leq t, P_1 \leq D_1)$$

$s \leq t$, and same for \tilde{F}_0 .

Joint distribution for P_j, D_j

Identifiability: Let

$$G_j = p(D_j)$$

and

$$V_j(s | t) = p(P_j \leq s, P_j < D_j | D_j = t).$$

$s < t$ (for the moment, ignoring regression on “ \mathbf{x} ”).

Under random censoring G_j and V_j are identifiable – just use the corresponding sample statistics.

Bivariate sub-distribution: together G_j & V_j define

$$\tilde{F}_1(s, t) = p(P_1 \leq s, D_1 \leq t, P_1 \leq D_1)$$

$s \leq t$, and same for \tilde{F}_0 .

Random prob measures, $F_1(s, t)$ & $F_0(s, t)$ imply \tilde{F}_1 & \tilde{F}_0 .

DDP mix of normals, as before

Copula $G(D_0, D_1)$

Copula: Link F_0 and F_1 with a normal copula.

Φ = standard normal c.d.f and

$\Phi_{2,\rho}$ = bivariate normal with correlation ρ .

Copula $G(D_0, D_1)$

Copula: Link F_0 and F_1 with a normal copula.

Φ = standard normal c.d.f and

$\Phi_{2,\rho}$ = bivariate normal with correlation ρ .

$$G(D_0, D_1; \rho) = \Phi_{2,\rho} [\Phi^{-1}\{G_0(D_0)\}, \Phi^{-1}\{G_1(D_1)\}]$$

Copula $G(D_0, D_1)$

Copula: Link F_0 and F_1 with a normal copula.

Φ = standard normal c.d.f and

$\Phi_{2,\rho}$ = bivariate normal with correlation ρ .

$$G(D_0, D_1; \rho) = \Phi_{2,\rho} [\Phi^{-1}\{G_0(D_0)\}, \Phi^{-1}\{G_1(D_1)\}]$$

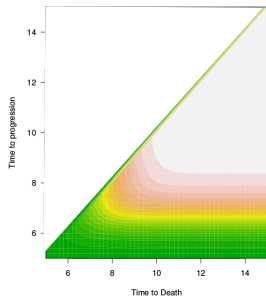
ρ is not identifiable – choice of ρ is an **assumption**.

Odds of progression

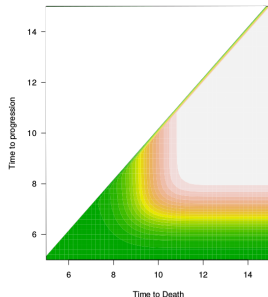
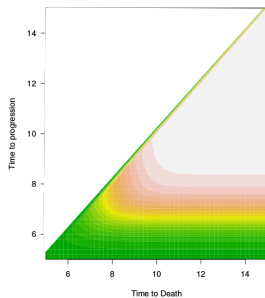
Then

$$\tau_{\mathbf{x}}(u) = \frac{\int_{P_1 < u} \int_{D_0 \geq u} \int_{D_1 \geq u} dV_1(P_1 | D_1, \mathbf{x}) dG_{\mathbf{x}}(D_0, D_1)}{\int_{P_0 < u} \int_{D_0 \geq u} \int_{D_1 \geq u} dV_0(P_0 | D_1), \mathbf{x} dG_{\mathbf{x}}(D_0, D_1)}$$

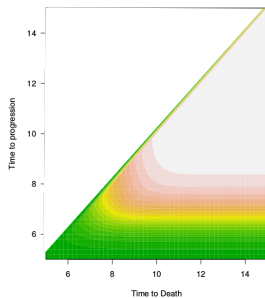
Results – Brain tumor study



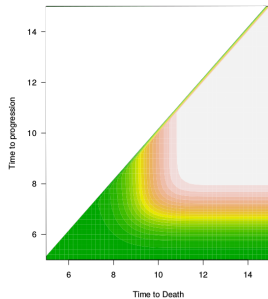
Results – Brain tumor study



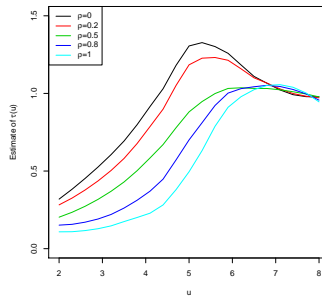
Results – Brain tumor study



$p(D_0, P_0)$



$p(D_1, P_1)$



$\tau(u)$ (on *log* scale!)

BNP regression by covariate-dependent partitions

Define BNP regression by

- 1 random partition, indexed by covariates;
- 2 cluster-specific sampling model.

→ next topic..

3. Classification

Categorical x_i : different subpopulations of interest

Aim: classify a new patient as $x_{n+1} = x \in \{0, 1\}$

Model:

$$y_i \mid x_i = 1 \sim F_x \text{ and } \{F_x; x = 0, 1\} \sim \text{DDP}$$

as before (GP simplifies to bivariate normal for $x \in \{0, 1\}$),
but ...

3. Classification

Categorical x_i : different subpopulations of interest

Aim: classify a new patient as $x_{n+1} = x \in \{0, 1\}$

Model:

$$y_i \mid x_i = 1 \sim F_x \text{ and } \{F_x; x = 0, 1\} \sim \text{DDP}$$

as before (GP simplifies to bivariate normal for $x \in \{0, 1\}$),
but ...

Simple augmentation: with

$$p(x_i = 1) = \pi$$

allows the desired ...

3. Classification

Categorical x_i : different subpopulations of interest

Aim: classify a new patient as $x_{n+1} = x \in \{0, 1\}$

Model:

$$y_i \mid x_i = 1 \sim F_x \text{ and } \{F_x; x = 0, 1\} \sim \text{DDP}$$

as before (GP simplifies to bivariate normal for $x \in \{0, 1\}$),
but ...

Simple augmentation: with

$$p(x_i = 1) = \pi$$

allows the desired ...

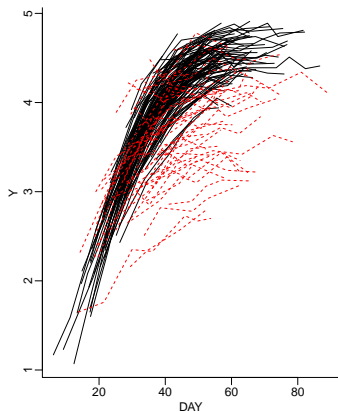
Classification: $p(x_{n+1} = 1 \mid \text{data})$ – that's all! (de la Cruz et al., 2007
ApplStat)

Example 3: Pregnancy classification

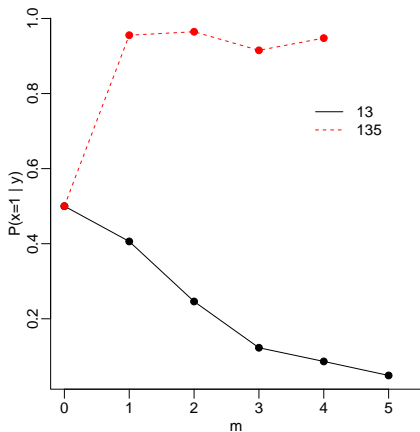
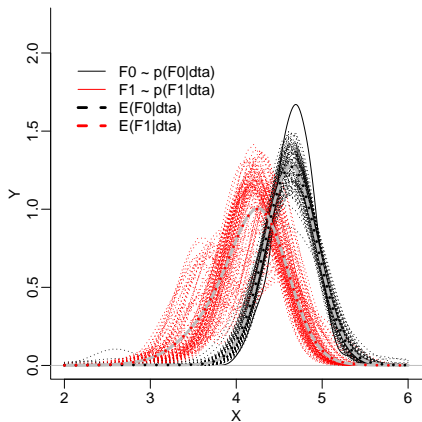
De la Cruz-Mesia et al. (2007, ApplStat)

Data: hormone data y_{ij} for $n = 173$ pregnant women, repeat mmt at times t_{ij} , $j = 1, \dots, n_i$

Subpopulations: $x_i = 0$, normal pregnancies, $n_0 = 124$ women
 $x_i = 1$, abnormal, $n_1 = 49$



Sampling model: $y_{ij} \mid x_i = x, \dots \sim N(m_{ij}, \sigma_x^2)$
with $m_{ij} = \theta_i / \{1 + e^{-(t_{ij} - \beta_{1x}) / \beta_{2x}}\}$



(a) $E(F_x | data)$

(b) $p(x_{n+1} = 1 | y_{n+1,1...m}, data)$

Estimated F_x under $x = 0$ (thick black curve) and $x = 1$ (thick red or grey) (panel a), and posterior probability $p(x_{n+1} = 1 | y_{n+1,1...m}, data)$

2. Clustering

Recall: DP Mixtures: convolution of discrete $F = \sum p_h \delta_{m_h}$ with (continuous) kernel, e.g., normal

$$G(y) = \int N(y | \theta, \sigma^2) dF(\theta), \quad F \sim \text{DP}$$

2. Clustering

Recall: DP Mixtures: convolution of discrete $F = \sum p_h \delta_{m_h}$ with (continuous) kernel, e.g., normal

$$\begin{aligned} G(y) &= \int N(y | \theta, \sigma^2) dF(\theta), \quad F \sim \text{DP} \\ &= \sum_{h=1}^{\infty} p_h N(y | m_h, \sigma) \end{aligned}$$

continuous $G(\cdot)$ (and hyperpar σ^2)

2. Clustering

Recall: DP Mixtures: convolution of discrete $F = \sum p_h \delta_{m_h}$ with (continuous) kernel, e.g., normal

$$\begin{aligned} G(y) &= \int N(y | \theta, \sigma^2) dF(\theta), \quad F \sim \text{DP} \\ &= \sum_{h=1}^{\infty} p_h N(y | m_h, \sigma) \end{aligned}$$

continuous $G(\cdot)$ (and hyperpar σ^2)

Latent vars: write $\int \dots dF(\theta)$ as hierarchical model

$$\begin{aligned} y_i | \theta_i &\sim N(\theta_i, \sigma^2), \quad i = 1, \dots, n \\ \theta_i | F &\sim F \end{aligned}$$

2. Clustering

Recall: DP Mixtures: convolution of discrete $F = \sum p_h \delta_{m_h}$ with (continuous) kernel, e.g., normal

$$\begin{aligned} G(y) &= \int N(y | \theta, \sigma^2) dF(\theta), \quad F \sim \text{DP} \\ &= \sum_{h=1}^{\infty} p_h N(y | m_h, \sigma) \end{aligned}$$

continuous $G(\cdot)$ (and hyperpar σ^2)

Latent vars: write $\int \dots dF(\theta)$ as hierarchical model

$$\begin{aligned} y_i | \theta_i &\sim N(\theta_i, \sigma^2), \quad i = 1, \dots, n \\ \theta_i | F &\sim F \end{aligned}$$

Notation: discrete $F \Rightarrow K \leq n$ unique θ_i 's = $\{\phi_1^*, \dots, \phi_K^*\}$.

2. Clustering

Recall: DP Mixtures: convolution of discrete $F = \sum p_h \delta_{m_h}$ with (continuous) kernel, e.g., normal

$$\begin{aligned} G(y) &= \int N(y | \theta, \sigma^2) dF(\theta), \quad F \sim \text{DP} \\ &= \sum_{h=1}^{\infty} p_h N(y | m_h, \sigma) \end{aligned}$$

continuous $G(\cdot)$ (and hyperpar σ^2)

Latent vars: write $\int \dots dF(\theta)$ as hierarchical model

$$\begin{aligned} y_i | \theta_i &\sim N(\theta_i, \sigma^2), \quad i = 1, \dots, n \\ \theta_i | F &\sim F \end{aligned}$$

Notation: discrete $F \Rightarrow K \leq n$ unique θ_i 's = $\{\phi_1^*, \dots, \phi_K^*\}$.

Latent indicators: $z_i = j$ iff $\theta_i = \phi_j^*$ match θ_i with ϕ_j^* 's.

2. Clustering

Recall: DP Mixtures: convolution of discrete $F = \sum p_h \delta_{m_h}$ with (continuous) kernel, e.g., normal

$$\begin{aligned} G(y) &= \int N(y | \theta, \sigma^2) dF(\theta), \quad F \sim \text{DP} \\ &= \sum_{h=1}^{\infty} p_h N(y | m_h, \sigma) \end{aligned}$$

continuous $G(\cdot)$ (and hyperpar σ^2)

Latent vars: write $\int \dots dF(\theta)$ as hierarchical model

$$\begin{aligned} y_i | \theta_i &\sim N(\theta_i, \sigma^2), \quad i = 1, \dots, n \\ \theta_i | F &\sim F \end{aligned}$$

Notation: discrete $F \Rightarrow K \leq n$ unique θ_i 's = $\{\phi_1^*, \dots, \phi_K^*\}$.

Latent indicators: $z_i = j$ iff $\theta_i = \phi_j^*$ match θ_i with ϕ_j^* 's.

Random Partition Models

Product partition model (PPM): cohesion functions $c(S_j)$ define similarity of a cluster,

$$p(\rho_n) \propto \prod_{j=1}^k c(S_j).$$

Hartigan (1990 Comm Stat), Barry and Hartigan (1993 JASA)

Random Partition Models

Product partition model (PPM): cohesion functions $c(S_j)$ define similarity of a cluster,

$$p(\rho_n) \propto \prod_{j=1}^k c(S_j).$$

Hartigan (1990 Comm Stat), Barry and Hartigan (1993 JASA)

Sampling model: conditional on partition ρ_n , assume exchangeability,

$$p(y^n \mid \rho, \phi^*) = \prod_{j=1}^k \left\{ \prod_{i \in S_j} p(y_i \mid \phi_j^*) \right\} \quad (*)$$

with cluster-specific parameters ϕ_j^*

Random Partition Models

Product partition model (PPM): cohesion functions $c(S_j)$ define similarity of a cluster,

$$p(\rho_n) \propto \prod_{j=1}^k c(S_j).$$

Hartigan (1990 Comm Stat), Barry and Hartigan (1993 JASA)

Sampling model: conditional on partition ρ_n , assume exchangeability,

$$p(y^n \mid \rho, \phi^*) = \prod_{j=1}^k \left\{ \prod_{i \in S_j} p(y_i \mid \phi_j^*) \right\} \quad (*)$$

with cluster-specific parameters ϕ_j^*

Prior $p(\phi_j^*)$: conjugate ...

Covariate-dependent PPM (PPM_x)

M et al. (2011 JCGS), Quintana et al. (2015 ScandJS)

Random partition: to favor clusters of patients with similar covariates,

Covariate-dependent PPM (PPM_x)

M et al. (2011 JCGS), Quintana et al. (2015 ScandJS)

Random partition: to favor clusters of patients with similar covariates, define $g(x_j^*) > 0$ to characterize the similarity of $\{x_i; i \in S_j\}$ with low values for bad clusters:

Covariate-dependent PPM (PPM_x)

M et al. (2011 JCGS), Quintana et al. (2015 ScandJS)

Random partition: to favor clusters of patients with similar covariates, define $g(x_j^*) > 0$ to characterize the similarity of $\{x_i; i \in S_j\}$ with low values for bad clusters:

$$p(\rho_n | x^n) \propto \prod_{j=1}^k g(x_j^*) \cdot c(S_j)$$

Covariate-dependent PPM (PPM_x)

M et al. (2011 JCGS), Quintana et al. (2015 ScandJS)

Random partition: to favor clusters of patients with similar covariates, define $g(x_j^*) > 0$ to characterize the similarity of $\{x_i; i \in S_j\}$ with low values for bad clusters:

$$p(\rho_n | x^n) \propto \prod_{j=1}^k g(x_j^*) \cdot c(S_j)$$

Similarity function: easy computation with

$$g(x_j^*) = \int \prod_{i \in S_j} q(x_i | \xi_j^*) q(\xi_j^*) d\xi_j^*$$

Covariate-dependent PPM (PPM_x)

M et al. (2011 JCGS), Quintana et al. (2015 ScandJS)

Random partition: to favor clusters of patients with similar covariates, define $g(x_j^*) > 0$ to characterize the similarity of $\{x_i; i \in S_j\}$ with low values for bad clusters:

$$p(\rho_n | x^n) \propto \prod_{j=1}^k g(x_j^*) \cdot c(S_j)$$

Similarity function: easy computation with

$$g(x_j^*) = \int \prod_{i \in S_j} q(x_i | \xi_j^*) q(\xi_j^*) d\xi_j^*$$

using, e.g., $q(x_i | \xi_i) = N(\xi_i^*, V)$ and $q(\xi_j^*) = N(\dots)$ for continuous x_i ,

Covariate-dependent PPM (PPM_x)

M et al. (2011 JCGS), Quintana et al. (2015 ScandJS)

Random partition: to favor clusters of patients with similar covariates, define $g(x_j^*) > 0$ to characterize the similarity of $\{x_i; i \in S_j\}$ with low values for bad clusters:

$$p(\rho_n | x^n) \propto \prod_{j=1}^k g(x_j^*) \cdot c(S_j)$$

Similarity function: easy computation with

$$g(x_j^*) = \int \prod_{i \in S_j} q(x_i | \xi_j^*) q(\xi_j^*) d\xi_j^*$$

using, e.g., $q(x_i | \xi_i) = N(\xi_i^*, V)$ and $q(\xi_j^*) = N(\dots)$ for continuous x_i , and similar conjugate choices for categorical, ordinal and counts.

Example

Example 4: Survival regression with PPMx

M, Quintana & Rosner (2011 JCGS) analyze data from a study (CALGB 9082) of breast cancer patients.

Treatment: high dose (A) versus low dose (B) chemotherapy

Example

Example 4: Survival regression with PPMx

M, Quintana & Rosner (2011 JCGS) analyze data from a study (CALGB 9082) of breast cancer patients.

Treatment: high dose (A) versus low dose (B) chemotherapy

Data: 765 patients randomized to A vs. B.

Example

Example 4: Survival regression with PPMx

M, Quintana & Rosner (2011 JCGS) analyze data from a study (CALGB 9082) of breast cancer patients.

Treatment: high dose (A) versus low dose (B) chemotherapy

Data: 765 patients randomized to A vs. B.

Response: time until progression or death

Example

Example 4: Survival regression with PPMx

M, Quintana & Rosner (2011 JCGS) analyze data from a study (CALGB 9082) of breast cancer patients.

Treatment: high dose (A) versus low dose (B) chemotherapy

Data: 765 patients randomized to A vs. B.

Response: time until progression or death

Covariates:

- *Categorical:* dose (A vs. B), menopausal status, estrogen use
- *Continuous:* age, initial tumor size,
- *Count:* number of positive lymph nodes

Example

Example 4: Survival regression with PPM_x

M, Quintana & Rosner (2011 JCGS) analyze data from a study (CALGB 9082) of breast cancer patients.

Treatment: high dose (A) versus low dose (B) chemotherapy

Data: 765 patients randomized to A vs. B.

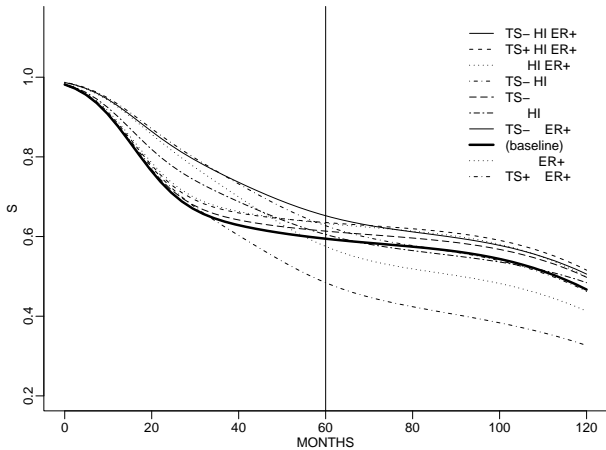
Response: time until progression or death

Covariates:

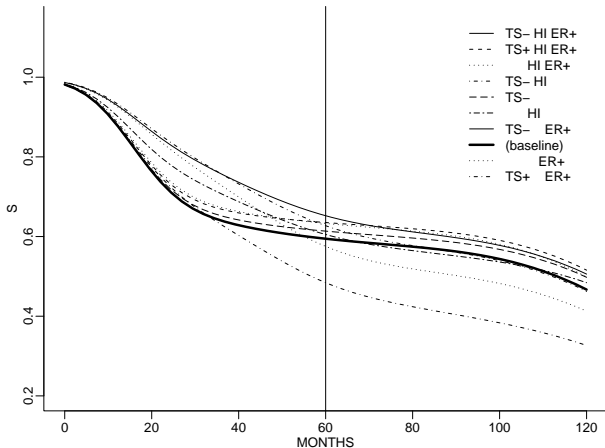
- *Categorical:* dose (A vs. B), menopausal status, estrogen use

- *Continuous:* age, initial tumor size,
- *Count:* number of positive lymph nodes

Model: PPM_x, with cluster-specific normal sampling



$S(t | x)$ by covariates



$S(t | x)$ by covariates

BNP regression: use the PPMx for BNP regression; allowing regression with variable dimension covariate vector!

Bayesian Subgroup analysis

Subgroup analysis problem: inference on exceptions from overall conclusion, typically for a clinical study, for

- a “benefitting population”,

vs.

- eligible population of the trial

Bayesian Subgroup analysis

Subgroup analysis problem: inference on exceptions from overall conclusion, typically for a clinical study, for

- a “benefitting population”,

vs.

- eligible population of the trial

Approaches :

- Treatment/cov interaction: Dixon and Simon (1991 Bmcs), Jones et al. (2011 ClinTrials)

Bayesian Subgroup analysis

Subgroup analysis problem: inference on exceptions from overall conclusion, typically for a clinical study, for

- a “benefitting population”,

vs.

- eligible population of the trial

Approaches :

- Treatment/cov interaction: Dixon and Simon (1991 Bmcs), Jones et al. (2011 ClinTrials)
- Tree based methods: Foster, Taylor & Ruberg (2011 StatMed)

Bayesian Subgroup analysis

Subgroup analysis problem: inference on exceptions from overall conclusion, typically for a clinical study, for

- a “benefitting population”,

vs.

- eligible population of the trial

Approaches :

- Treatment/cov interaction: Dixon and Simon (1991 Bmcs), Jones et al. (2011 ClinTrials)
- Tree based methods: Foster, Taylor & Ruberg (2011 StatMed)
- Model selection: Berger, Wang and Shen (2014, J Biopharm Stat), Sivaganesan et al. (2011 StatMed)

Bayesian Subgroup analysis

Subgroup analysis problem: inference on exceptions from overall conclusion, typically for a clinical study, for

- a “benefitting population”,

vs.

- eligible population of the trial

Approaches :

- Treatment/cov interaction: Dixon and Simon (1991 Bmcs), Jones et al. (2011 ClinTrials)
- Tree based methods: Foster, Taylor & Ruberg (2011 StatMed)
- Model selection: Berger, Wang and Shen (2014, J Biopharm Stat), Sivaganesan et al. (2011 StatMed)
- Decision problem: next slides...

Decision Problem

Data: response y_i , covariates $x_i = (x_{i1}, \dots, x_{ip})$.

Decision Problem

Data: response y_i , covariates $x_i = (x_{i1}, \dots, x_{ip})$.

Actions: Report a subgroup of patients who most benefit from the experimental therapy:

$$\mathbf{a} = (I, \mathbf{x}^*),$$

Covariates: $I \subset \{1, \dots, p\}$

Levels: $\mathbf{x}^* = (x_j^*, j \in I)$,

(possibly restrict continuous x_j^* to fixed thresholds)

Decision Problem

Data: response y_i , covariates $x_i = (x_{i1}, \dots, x_{ip})$.

Actions: Report a subgroup of patients who most benefit from the experimental therapy:

$$\mathbf{a} = (I, \mathbf{x}^*),$$

Covariates: $I \subset \{1, \dots, p\}$

Levels: $\mathbf{x}^* = (x_j^*, j \in I)$,

(possibly restrict continuous x_j^* to fixed thresholds)

Decision problem: separate **inference** (predicting y_{n+1}), with flexible model

Decision Problem

Data: response y_i , covariates $x_i = (x_{i1}, \dots, x_{ip})$.

Actions: Report a subgroup of patients who most benefit from the experimental therapy:

$$\mathbf{a} = (I, \mathbf{x}^*),$$

Covariates: $I \subset \{1, \dots, p\}$

Levels: $\mathbf{x}^* = (x_j^*, j \in I)$,

(possibly restrict continuous x_j^* to fixed thresholds)

Decision problem: separate **inference** (predicting y_{n+1}), with flexible model
vs.

decision (report subpopulation), parsimoniously

Decision Problem

Data: response y_i , covariates $x_i = (x_{i1}, \dots, x_{ip})$.

Actions: Report a subgroup of patients who most benefit from the experimental therapy:

$$\mathbf{a} = (I, \mathbf{x}^*),$$

Covariates: $I \subset \{1, \dots, p\}$

Levels: $\mathbf{x}^* = (x_j^*, j \in I)$,

(possibly restrict continuous x_j^* to fixed thresholds)

Decision problem: separate **inference** (predicting y_{n+1}), with flexible model
vs.

decision (report subpopulation), parsimoniously

- no need for multiplicity control
- arbitrary prob model
- disentangle stat significance vs. clinical relevance
- allow for variable # covs.

Utility: we favor a subpopulation with difference (relative to the overall population) in **trt effect**, large **size** and parsimonious description with **few covariates**.

- *Event time:* e.g., for an $y_i = \text{PFS}$ (event time), this could be based on log hazard ratio

Utility: we favor a subpopulation with difference (relative to the overall population) in **trt effect**, large **size** and parsimonious description with **few covariates**.

- *Event time:* e.g., for an $y_i = \text{PFS}$ (event time), this could be based on log hazard ratio

$$u(\mathbf{a}, \theta) = (\text{LR}(\mathbf{a}, \theta) - \beta) \cdot \frac{n(\mathbf{a})^\alpha}{(|I| + 1)^\gamma} \quad (1)$$

where θ are parameters that index the sampling model.

Utility: we favor a subpopulation with difference (relative to the overall population) in **trt effect**, large **size** and parsimonious description with **few covariates**.

- *Event time:* e.g., for an $y_i = \text{PFS}$ (event time), this could be based on log hazard ratio

$$u(a, \theta) = (\text{LR}(a, \theta) - \beta) \cdot \frac{n(\mathbf{a})^\alpha}{(|I| + 1)^\gamma} \quad (1)$$

where θ are parameters that index the sampling model.

- *Continuous outcome:* e.g., % tumor shrinkage, this could be based on predictive average treatment effect (PATE),

Utility: we favor a subpopulation with difference (relative to the overall population) in **trt effect**, large **size** and parsimonious description with **few covariates**.

- *Event time:* e.g., for an $y_i = \text{PFS}$ (event time), this could be based on log hazard ratio

$$u(a, \theta) = (\text{LR}(a, \theta) - \beta) \cdot \frac{n(\mathbf{a})^\alpha}{(|I| + 1)^\gamma} \quad (1)$$

where θ are parameters that index the sampling model.

- *Continuous outcome:* e.g., % tumor shrinkage, this could be based on predictive average treatment effect (PATE), averaged over x_i and already averaged w.r.t. $p(\theta | \text{data})$.

Utility: we favor a subpopulation with difference (relative to the overall population) in **trt effect**, large **size** and parsimonious description with **few covariates**.

- *Event time:* e.g., for an $y_i = \text{PFS}$ (event time), this could be based on log hazard ratio

$$u(a, \theta) = (\text{LR}(a, \theta) - \beta) \cdot \frac{n(a)^\alpha}{(|I| + 1)^\gamma} \quad (1)$$

where θ are parameters that index the sampling model.

- *Continuous outcome:* e.g., % tumor shrinkage, this could be based on predictive average treatment effect (PATE), averaged over x_i and already averaged w.r.t. $p(\theta | \text{data})$.

$$U(a) = \begin{cases} \{\text{PATE}_{SS}(a) - \beta\} \cdot \frac{|n(a)+1|^\alpha}{(|I|+1)^\gamma} & \text{if } a \neq H_0 \\ u_0 & \text{if } a = H_0, \end{cases}$$

where H_0, H_1 are special actions,

Utility: we favor a subpopulation with difference (relative to the overall population) in **trt effect**, large **size** and parsimonious description with **few covariates**.

- *Event time:* e.g., for an $y_i = \text{PFS}$ (event time), this could be based on log hazard ratio

$$u(a, \theta) = (\text{LR}(a, \theta) - \beta) \cdot \frac{n(a)^\alpha}{(|I| + 1)^\gamma} \quad (1)$$

where θ are parameters that index the sampling model.

- *Continuous outcome:* e.g., % tumor shrinkage, this could be based on predictive average treatment effect (PATE), averaged over x_i and already averaged w.r.t. $p(\theta | \text{data})$.

$$U(a) = \begin{cases} \{\text{PATE}_{SS}(a) - \beta\} \cdot \frac{|n(a)+1|^\alpha}{(|I|+1)^\gamma} & \text{if } a \neq H_0 \\ u_0 & \text{if } a = H_0, \end{cases}$$

where H_0, H_1 are special actions,
with $\beta > 0$ a fixed clinically decided threshold and $n(a)$ is the size of the subpopulation.

θ indexes the sampling model

Utility: we favor a subpopulation with difference (relative to the overall population) in **trt effect**, large **size** and parsimonious description with **few covariates**.

- *Event time:* e.g., for an $y_i = \text{PFS}$ (event time), this could be based on log hazard ratio

$$u(a, \theta) = (\text{LR}(a, \theta) - \beta) \cdot \frac{n(a)^\alpha}{(|I| + 1)^\gamma} \quad (1)$$

where θ are parameters that index the sampling model.

- *Continuous outcome:* e.g., % tumor shrinkage, this could be based on predictive average treatment effect (PATE), averaged over x_i and already averaged w.r.t. $p(\theta | \text{data})$.

$$U(a) = \begin{cases} \{\text{PATE}_{SS}(a) - \beta\} \cdot \frac{|n(a)+1|^\alpha}{(|I|+1)^\gamma} & \text{if } a \neq H_0 \\ u_0 & \text{if } a = H_0, \end{cases}$$

where H_0, H_1 are special actions, with $\beta > 0$ a fixed clinically decided threshold and $n(a)$ is the size of the subpopulation.

θ indexes the sampling model (**any** model for $p(y | x, \theta)$)

Bayes rule: Report $a^* = \arg \max_a \int u(a, \theta) dp(\theta | data)$

Bayes rule: Report $a^* = \arg \max_a \int u(a, \theta) dp(\theta | data)$

Alternative utility: Foster, Taylor & Ruberg (2011, StatMed) use

$Q(A) = \text{enhanced treatment effect} - \text{average trt effect}$

and sensitivity and specificity to evaluate a reported subpopulation A .

Bayes rule: Report $a^* = \arg \max_a \int u(a, \theta) dp(\theta | data)$

Alternative utility: Foster, Taylor & Ruberg (2011, StatMed) use

$$Q(A) = \text{enhanced treatment effect} - \text{average trt effect}$$

and sensitivity and specificity to evaluate a reported subpopulation A .

Model: Decision problem and solution meaningful for **any** model.

3. Probability Model

Flexible BNP model. The BNP model “is always right.”

- *Event time*: for example, PPM_x for the event time
- *Continuous outcome*: e.g., DDP, BART

Example 5: Phase III Study of NSCL Patients

Morita & M, 2017 Bmcs

Patients: advanced non-small cell lung cancer, $n = 267$

Example 5: Phase III Study of NSCL Patients

Morita & M, 2017 Bmcs

Patients: advanced non-small cell lung cancer, $n = 267$

Treatment: carboplatin (N) ($n_0 = 130$) vs. paclitaxel + carboplatin (C) ($n_1 = 137$).

Example 5: Phase III Study of NSCL Patients

Morita & M, 2017 Bmcs

Patients: advanced non-small cell lung cancer, $n = 267$

Treatment: carboplatin (N) ($n_0 = 130$) vs. paclitaxel + carboplatin (C) ($n_1 = 137$).

Baseline covariates: pharmacologically relevant gene expressions, including 16 mRNA (mR1 - mR16) and 1 protein (Pn1) expression levels ($p = 17$).

Example 5: Phase III Study of NSCL Patients

Morita & M, 2017 Bmcs

Patients: advanced non-small cell lung cancer, $n = 267$

Treatment: carboplatin (N) ($n_0 = 130$) vs. paclitaxel + carboplatin (C) ($n_1 = 137$).

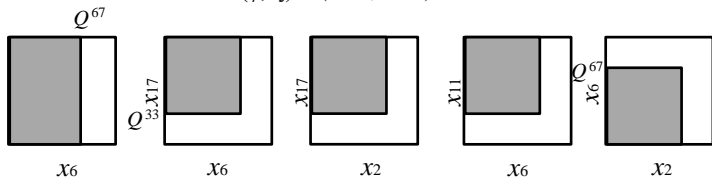
Baseline covariates: pharmacologically relevant gene expressions, including 16 mRNA (mR1 - mR16) and 1 protein (Pn1) expression levels ($p = 17$).

Outcome: $y_i = \max TS\%$ (max tumor size shrinkage from baseline)

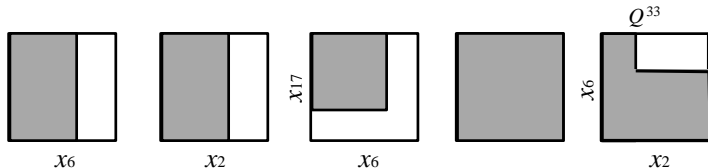
Results

Implement subgroup analysis for the phase III NSCL trial, restricting subgroups to $|I| \leq 2$ covariates.

$$(\phi, \zeta) = (0.25, 0.15)$$



$$(\phi, \zeta) = (0.35, 0.25)$$



Example 6: A basket trial design for targeted therapies

Xu et al. (2018 Biometrical J)

Subgroup analysis with a purpose.

IMPACT II: patients across different cancers. Based on molecular alterations patients are eligible for certain targeted therapies (TT)

Subgroup analysis: find subgroup of tumor/mutation pairs who most benefit from TT

Selecting the subpopulations

- Based on a flexible probability model: PPM_x

Selecting the subpopulations

- Based on a flexible probability model: PPM_x
- Utility function: $u(a, \dots)$ (1) for event time, PFS

Selecting the subpopulations

- Based on a flexible probability model: PPMx
- Utility function: $u(a, \dots)$ (1) for event time, PFS
- Report the subpopulations with largest expected utility
- Adaptive treatment allocation

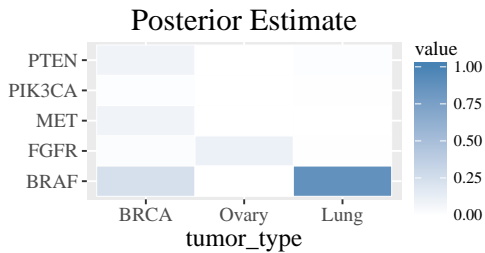
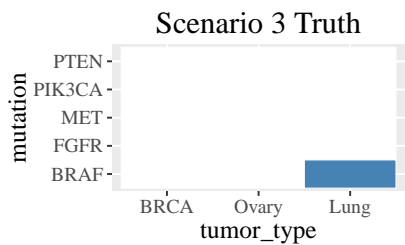
Simulation

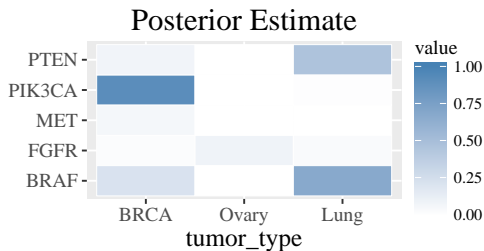
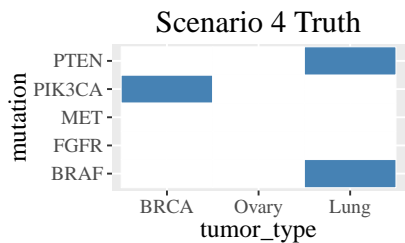
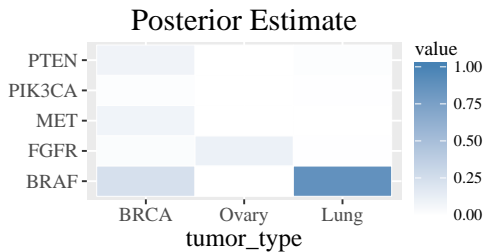
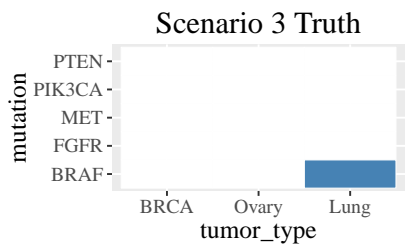
6 scenarios: overall treatment effect (trt);
interaction $z \times \text{mutation} \times \text{tumor}$,
 $z \in \{0, 1\}$, mutation $\in \{\text{BRAF}, \text{PIK3CA}, \text{PTEN}\}$,
tumors $\in \{\text{BRCA}, \text{Lung}, \text{Ovary}\}$.

Simulation

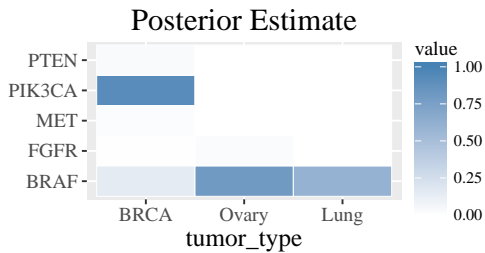
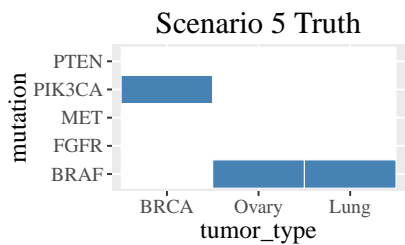
6 scenarios: overall treatment effect (trt);
interaction $z \times \text{mutation} \times \text{tumor}$,
 $z \in \{0, 1\}$, mutation $\in \{\text{BRAF}, \text{PIK3CA}, \text{PTEN}\}$,
tumors $\in \{\text{BRCA}, \text{Lung}, \text{Ovary}\}$.

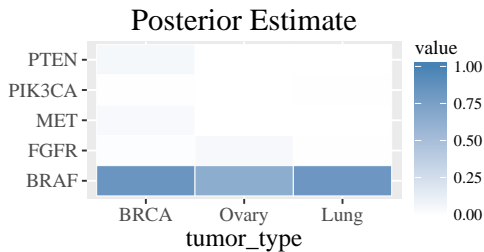
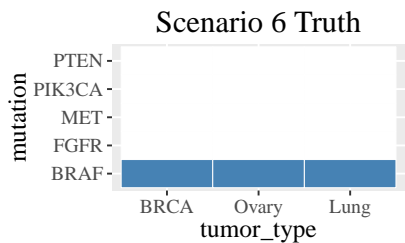
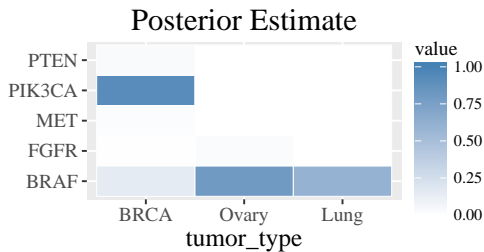
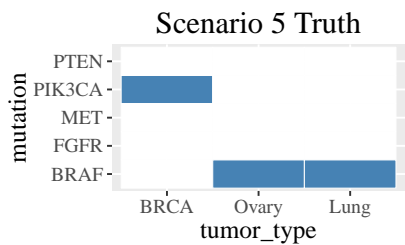
| | trt | Interactions (coefficient) |
|-------|-----|---|
| H_0 | 0 | none |
| H_1 | 0.4 | none |
| 3 | 0 | BRAF*Lung*z (0.4) |
| 4 | 0 | PIK3CA*BRCA*z (0.3), BRAF*Lung*z (0.3) PTEN*Lung*z(0.4) |
| 5 | 0 | PIK3CA*BRCA*z (0.3), BRAF*Ovary*z (0.4) BRAF*Lung*z(0.3) |
| 6 | 0 | BRAF*BRCA(0.4), BRAF*Ovary*z (0.3), BRAF*Lung*z(0.4) |





left = truth; right = estimate as $p(a)$ over repeat sim.





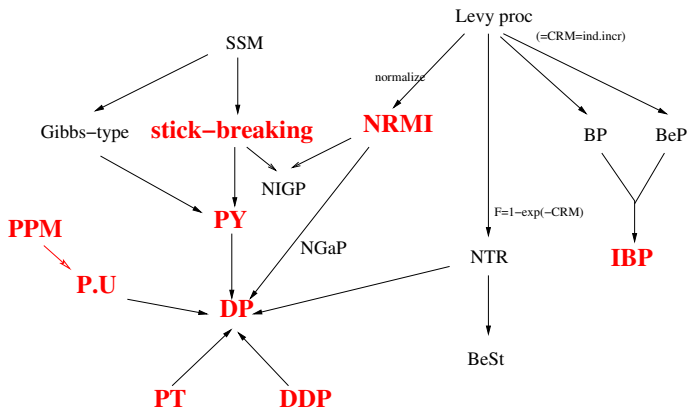
left = truth; right = estimate as $p(a)$ over repeat sim.

Summary

- **Definition:** BNP = prob models for infinite dim parameters.

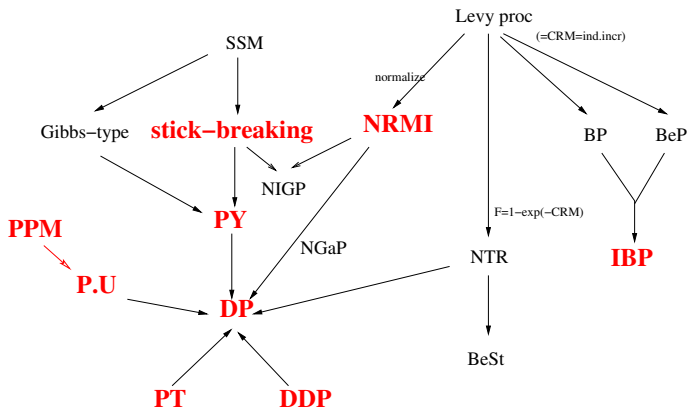
Summary

- **Definition:** BNP = prob models for infinite dim parameters.



Summary

- **Definition:** BNP = prob models for infinite dim parameters.



- Flexible models for full probabilistic description of all uncertainties
- Computation intensive; nonsense in – rubbish out :-)