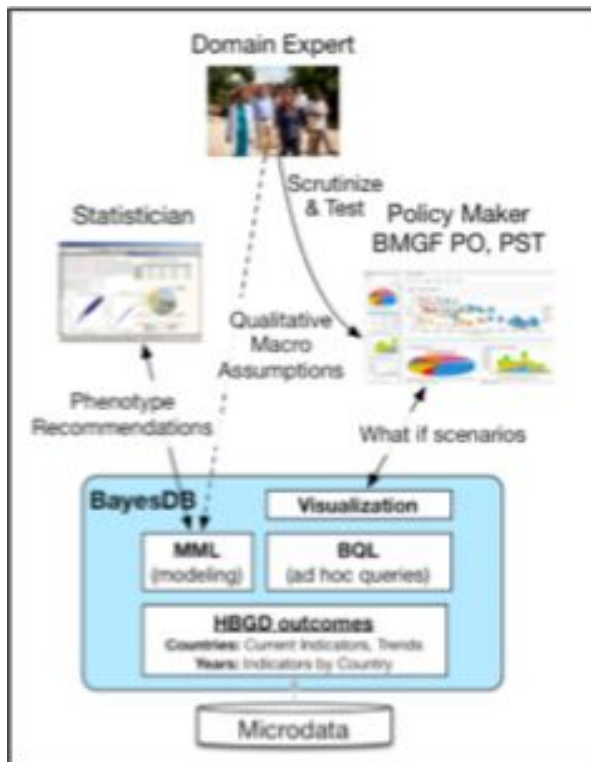


Probabilistic Computing Lab

Belhal KARIMI



The need for augmented intelligence



Policy advocate

"What are the comparable countries to Kenya in terms of everything we know about the malnutrition rate of infants?"

Domain expert

"Recent work in development economics suggests sanitation standards influence growth stunting in India but not in Africa."

Field researcher

"Here is new data on ~10,000 children in Bangladesh. Please update all relevant models and inform stakeholders."

Statistician

"Despite what the economists think, the p-value for this hypothesis test indicates that my mixed-effects model's finding of two different country clusters with respect to longitudinal variation in sanitation standards is not actually significant."

Main Research: Convergence accuracy of an inference program

- Inference program
- Which algorithm?
- Measure of accuracy:

- KL

$$D_{\text{KL}}(Q_a(X) || P(X|y)) + \mathcal{L}(Q_a) = \log P(y)$$

- ELBO

$$\mathcal{L}(Q_a) = \mathbb{E}_{x \sim Q_a} \left[\log \frac{\tilde{P}(x|y)}{Q_a(x)} \right]$$

Even if $\log P$ is not known, this holds:

Algorithm 1 Single-particle MCMC algorithm

Require: K_1, \dots, K_{T-1} are MCMC kernels and $Q_1^{aug}, \dots, Q_T^{aug}$ are augmentation distributions compatible with problem instance defined by $\tilde{P}(\theta, z_{1:T}, y_{1:T})$

```

1:  $w \leftarrow 1$ 
2:  $\theta^{(0)} \sim P(\Theta)$  ▷ Sample  $\theta$  from the prior
3:  $z_1^{(0)} \sim Q_1^{aug}(z_1 | \theta^{(0)}, y_1)$ 
4: for  $t \leftarrow 1$  to  $T - 1$  do
5:    $(\theta^{(t)}, z_{1:t}^{(t)}) \sim K_t(\theta, z_{1:t} | \theta^{(t-1)}, z_{1:t}^{(t-1)})$  ▷ Run the MCMC kernel  $K_t$ 
6:    $z_{t+1}^{(t)} \sim Q_{t+1}^{aug}(z_{t+1} | \theta^{(t)}, z_{1:t}^{(t)}, y_{1:t+1})$  ▷ Augment the state space
7:    $w \leftarrow w \cdot \frac{P(z_{t+1}^{(t)}, y_{t+1} | \theta^{(t)}, z_{1:t}^{(t)}, y_{1:t})}{Q_{t+1}^{aug}(z_{t+1}^{(t)} | \theta^{(t)}, z_{1:t}^{(t)}, y_{1:t+1})}$ 
8: end for
9: return  $(\theta^{(T-1)}, z_{1:T}^{(T-1)}), w$  ▷ Return the sample  $(\theta^{(T-1)}, z_{1:T}^{(T-1)}) \sim Q_a$  and the weight  $w$ 

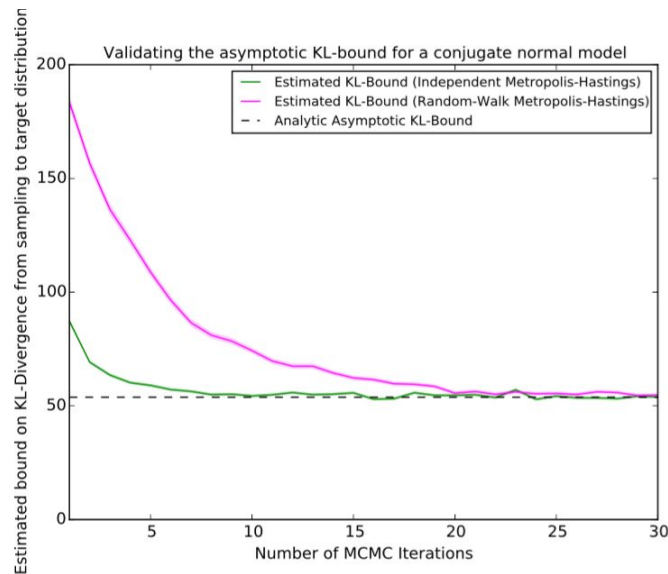
```

$$\mathcal{L}(Q_{a_1}) \geq \mathcal{L}(Q_{a_2}) \implies D_{\text{KL}}(Q_{a_1}(X) || P(X|y)) \leq D_{\text{KL}}(Q_{a_2}(X) || P(X|y))$$

Main Research: Convergence accuracy of an inference program

- Two proofs
 - KL does not increase by applying Kernel
 - Log weight is a good estimate of lower bound on ELBO

- Result
 - Independent Metropolis Hastings
 - Random walk Metropolis Hastings



Gates Foundation Case

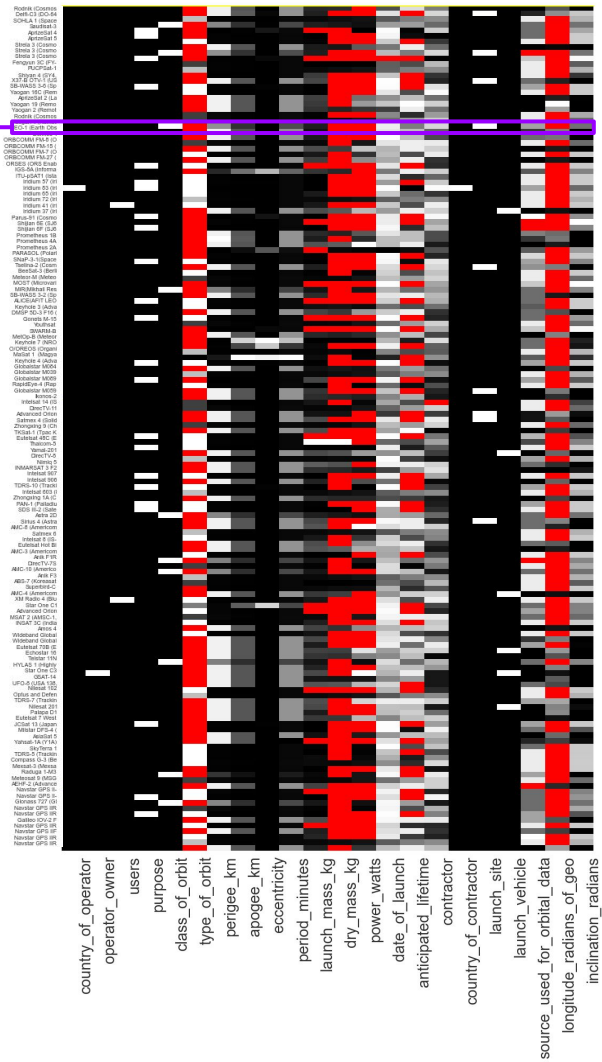
- Setting: UCS Satellites Database
 - 1167 rows (satellites) and 23 columns
 - Illustrations using 150 row subsample
 - Variables include, electrical, geopolitical, kinematic characteristics
 - Engineering note:
 - Schematics come from cleaned 'lovecat' states
 - Predictions come from 'gpmcc' states
- BayesDB capabilities illustrated:
 - Representing high-dimensional, incomplete, heterogeneously typed data
 - Estimating pairwise dependence probabilities from multiple GPMS
 - Generating simulations conditioned on hypotheticals

UCS Satellites Database: Raw Data

Data for Compass M4

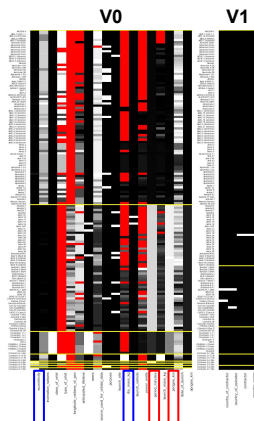
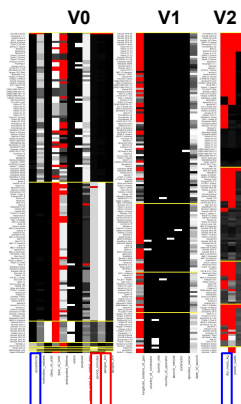
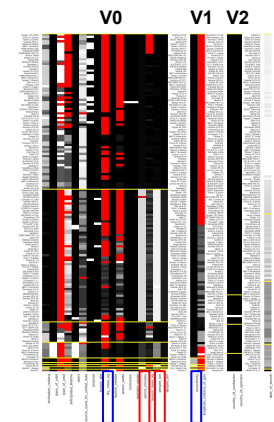
	0
Name	Compass M4 (Beidou 2-13)
Country_of_Operator	China (PR)
Operator_Owner	Chinese Defense Ministry
Users	Military
Purpose	Navigation/Global Positioning
Class_of_Orbit	MEO
Type_of_Orbit	NaN
Perigee_km	21452
Apogee_km	21603
Eccentricity	0.00271
Period_minutes	773.21
Launch_Mass_kg	2200
Dry_Mass_kg	NaN
Power_watts	NaN
Date_of_Launch	41027
Anticipated_Lifetime	8
Contractor	Space Technology Research Institute (part of C...
Country_of_Contractor	China (PR)
Launch_Site	Xichang Satellite Launch Center
Launch_Vehicle	Long March 3B
Source_Used_for_Orbital_Data	ZARYA
longitude_radians_of_geo	NaN
Inclination_radians	0.961676

Red are nans



Variable	Type
Country_of_Operator	categorical
Operator_Owner	categorical
Users	categorical
Purpose	categorical
Class_of_Orbit	categorical
Type_of_Orbit	categorical
Perigee_km	normal
Apogee_km	normal
Eccentricity	normal
Period_minutes	normal
Launch_Mass_kg	normal
Dry_Mass_kg	normal
Power_watts	normal
Date_of_Launch	normal
Anticipated_Lifetime	normal
Contractor	categorical
Country_of_Contractor	categorical
Launch_Site	categorical
Launch_Vehicle	categorical
Source_Used_for_Orbital_Data	categorical
longitude_radians_of_geo	normal
Inclination_radians	normal

UCS Satellites Database: Relation between Dependence Probability Heatmap and clustering

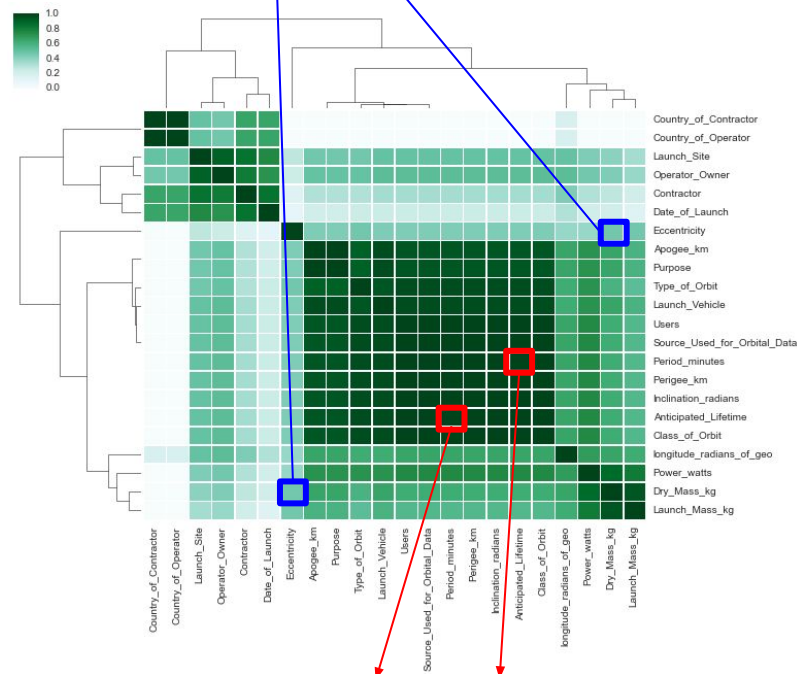


— Anticipated_lifetime
and Period_minutes

— Eccentricity and Dry_mass

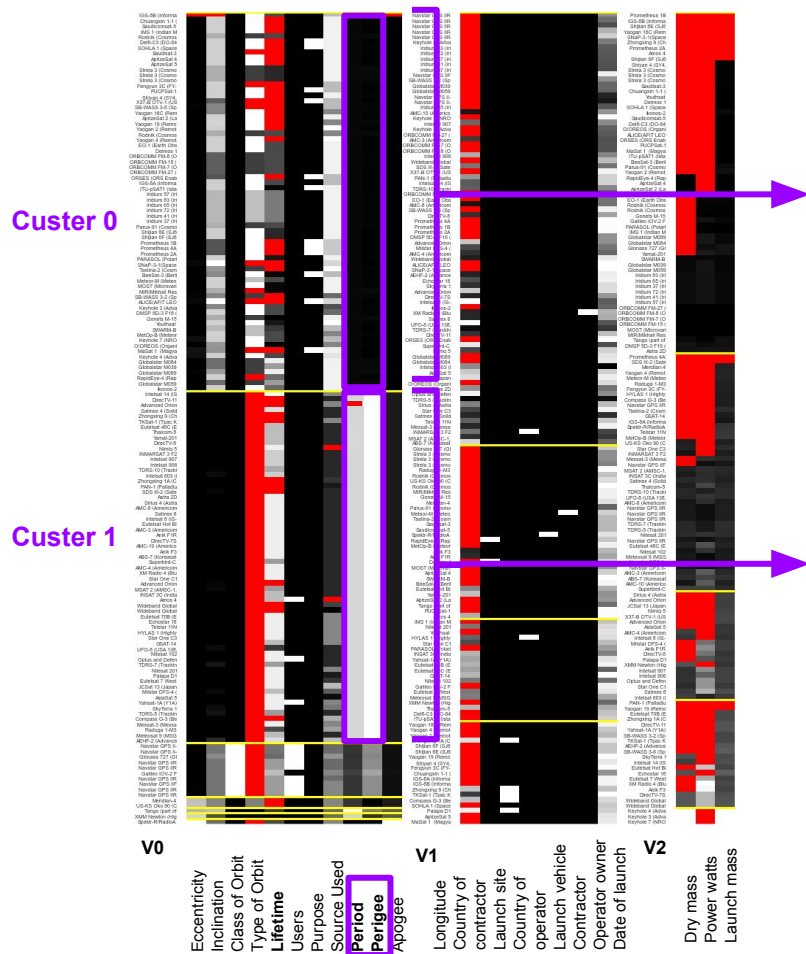
ESTIMATE DEPENDENCE PROBABILITY
FROM PAIRWISE COLUMNS OF generator

$$P(\text{eccentricity} \not\perp \text{dry_mass}) = 1/3$$



$$P(\text{Anticipated_lifetime} \not\perp \text{Period_minutes}) = 3/3$$

UCS Satellites Database: Generating simulations conditioned on hypotheticals

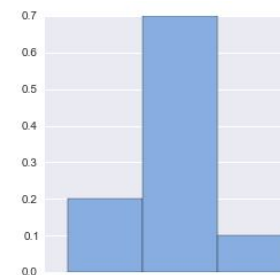


$p(\text{cluster}|\text{lifetime} = 3)$



Cluster 0 Cluster 1 Others

$p(\text{cluster}|\text{lifetime} = 15)$



Cluster 0 Cluster 1 Others

SIMULATE **Period_minutes**,
Perigee_km FROM generator GIVEN
Anticipated_Lifetime=3



$\sim p(\text{period, perigee}|\text{lifetime}=3)$

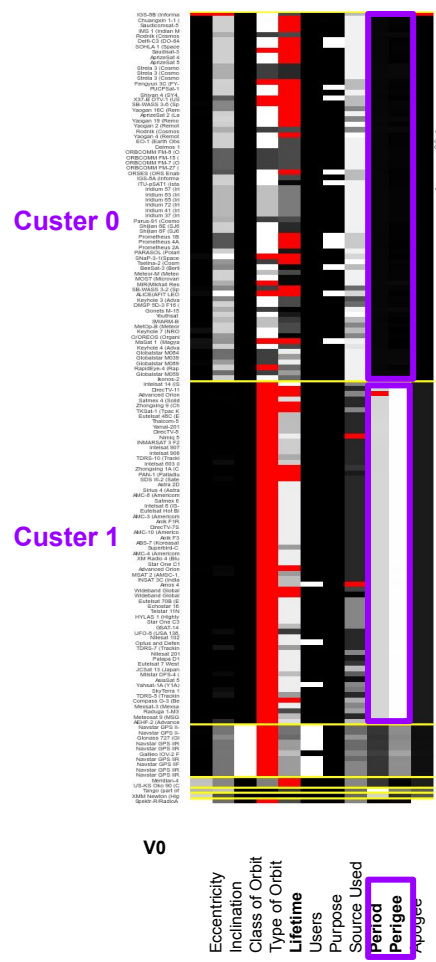
SIMULATE **Period_minutes**,
Perigee_km FROM generator GIVEN
Anticipated_Lifetime=15



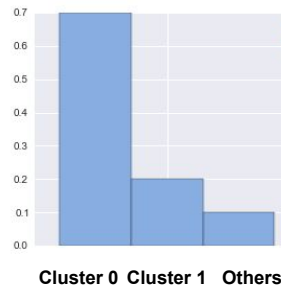
$\sim p(\text{period, perigee}|\text{lifetime}=15)$

UCS Satellites Database:

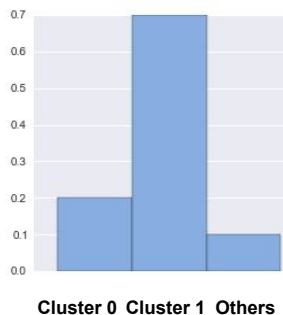
Posterior distribution vs. CC clustering



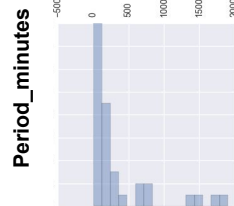
$p(\text{cluster}|\text{lifetime} = 3)$



$p(\text{cluster}|\text{lifetime} = 15)$



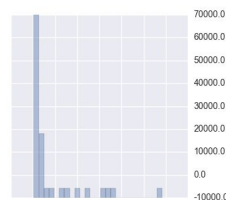
Period_minutes



Period_minutes

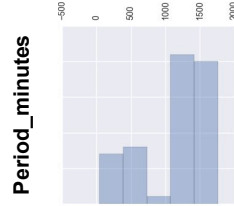
SIMULATE
Period_minutes,
Perigee_km FROM
generator GIVEN
Anticipated_Lifetime=3

Predictive checking



Perigee_km

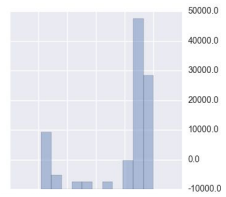
Period_minutes



Period_minutes

SIMULATE
Period_minutes,
Perigee_km FROM
generator GIVEN
Anticipated_Lifetime=15

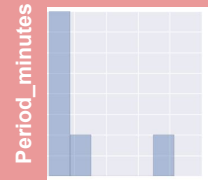
Predictive checking



Perigee_km

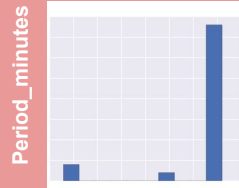
Real data

Period_minutes



SELECT Period_minutes,
Perigee_km FROM table
WHERE
Anticipated_Lifetime=3

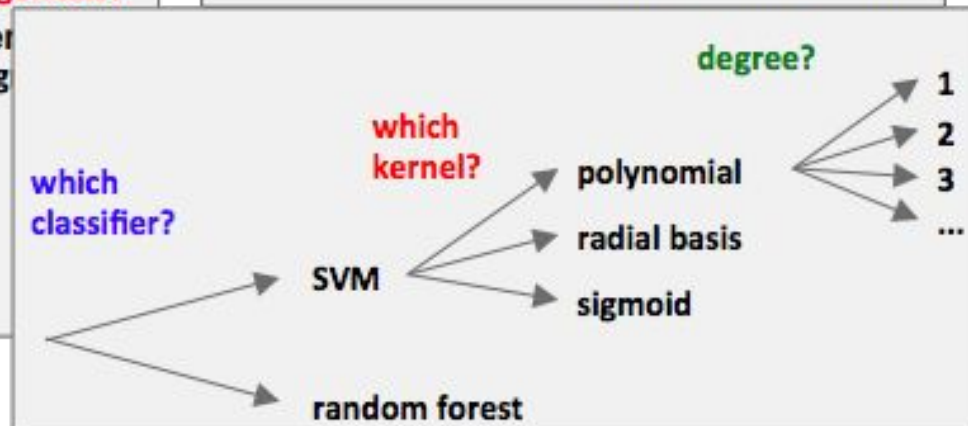
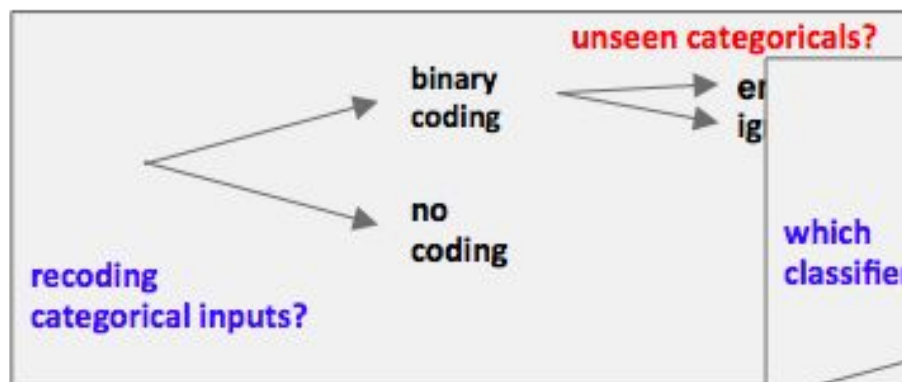
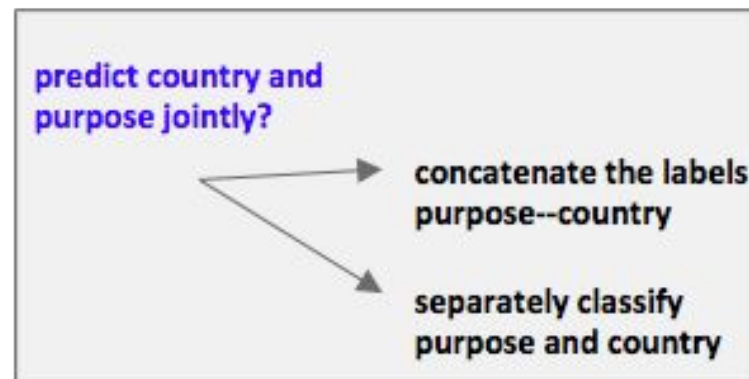
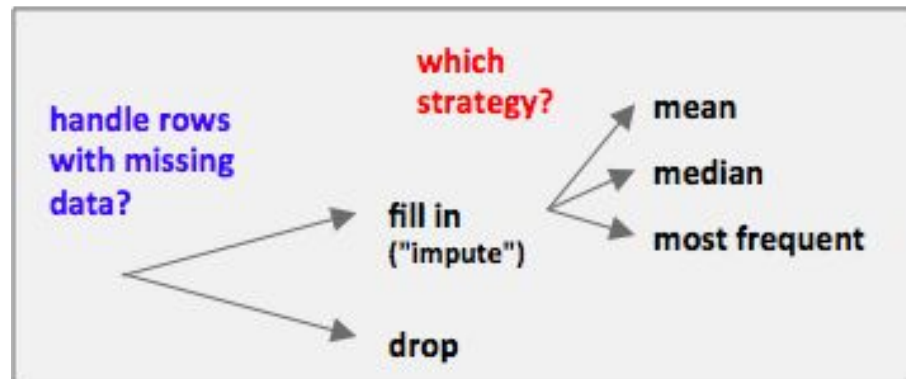
Period_minutes



SELECT Period_minutes,
Perigee_km FROM table
WHERE
Anticipated_Lifetime=15

BACKUP

Machine learning requires many decisions



Machine learning results are unstable

Approach 1

drop missing, no coding, random forest, separate classifiers

Simulations	Frequency
Egypt-Earth Science	9
Egypt-Earth/Space Science	5
Egypt-Astrophysics/Earth Science	3
Canada-Earth Science	3

Probably Egypt, definitely science

Approach 2

impute missing, binary coding, svm, joint classification

No idea

Approach 3

impute missing, no coding, random forest, separate classifiers

India-Meteorology	11
ESA-Meteorology	3
India-Communications	3
India-Earth Science	1
China (PR)-Space Physics	1
Russia-Space Physics	1

Probably India, probably science

STATISTICIAN

"Use the data from this .CSV file."

"Choose whatever data types you think are reasonable --- I don't have any knowledge about that."

"Build me a quick-and-dirty ensemble of models that gives me some ability to quantify uncertainty."

MML

```
CREATE POPULATION satellites  
FROM ucs_satellites.csv
```

```
CREATE METAMODEL ON satellites  
USING default_metamodel( GUESS(*) );
```

```
INITIALIZE 16 GENERATIVE POPULATION MODELS  
FOR satellites;  
ANALYZE satellites FOR 4 MINUTES;
```

Log weight is an estimate of ELBO

Since first proof

$$\begin{aligned} D_{\text{KL}}(\mathbf{Q}_K^{mcmc}(X^{(T)}) || P_{y_{1:T}}) &\leq D_{\text{KL}}(\mathbf{Q}_K^{mcmc}(X^{(T-1)}) || P_{y_{1:T}}) \\ \log p(y) - \mathcal{L}_{P_y}(\mathbf{Q}_K^{mcmc}(X^{(T)})) &\leq \log p(y) - \mathcal{L}_{P_y}(\mathbf{Q}_K^{mcmc}(X^{(T-1)})) \\ \mathcal{L}_{P_y}(\mathbf{Q}_K^{mcmc}(X^{(T-1)})) &\leq \mathcal{L}_{P_y}(\mathbf{Q}_K^{mcmc}(X^{(T)})) \end{aligned}$$

Note that $\mathbf{Q}_K^{mcmc}(X^{(T-1)})$ and $P_y(X^{(T-1)})$ are the marginal distributions of $X^{(T-1)}$ for \mathbf{Q}_K^{mcmc} and $\mathbf{P}_{y,K}^{mcmc}$, respectively. Therefore: $D_{\text{KL}}(\mathbf{Q}_K^{mcmc}(X^{(T-1)}) || P_y) \leq D_{\text{KL}}(\mathbf{Q}_K^{mcmc} || \mathbf{P}_{y,K}^{mcmc})$

So

$$\begin{aligned} \log p(y) - \mathcal{L}_{P_y}(\mathbf{Q}_K^{smc}(X^{(T-1)})) &\leq \log p(y) - \mathcal{L}_{\mathbf{P}_{y,K}^{smc}}(\mathbf{Q}_K^{smc}) \\ \mathcal{L}_{\mathbf{P}_{y,K}^{smc}}(\mathbf{Q}_K^{smc}) &\leq \mathcal{L}_{P_y}(\mathbf{Q}_K^{smc}(X^{(T-1)})) \end{aligned}$$

and finally

$$\begin{aligned} \mathbb{E}_{x^{(0:T-1)} \sim \mathbf{Q}_K^{smc}} \left[\log \frac{\tilde{\mathbf{p}}_{y,K}^{smc}(x^{(0:T-1)})}{\mathbf{q}_K^{smc}(x^{(0:T-1)})} \right] &\leq \mathcal{L}_{P_y}(\mathbf{Q}_K^{smc}(X^{(T-1)})) \\ \mathbb{E} [\log w] &\leq \mathcal{L}_{P_y}(\mathbf{Q}_K^{smc}(X^{(T-1)})) \end{aligned}$$

$\mathbb{E} [\log w] \leq \mathcal{L}_{P_y}(\mathbf{Q}_K^{smc}(X^{(T-1)})) \leq \mathcal{L}_{P_y}(\mathbf{Q}_K^{smc}(X^{(T)})) = \mathcal{L}_{P_y}(Q_K^{smc})$