# BayesDB: A probabilistic programming system for querying the probable implications of data

**Vikash Mansinghka**                              VKM@MIT.EDU
**Richard Tibbetts**                               TIBBETTS@MIT.EDU
**Jay Baxter**                                     JBAXTER@MIT.EDU
*Computer Science & Artificial Intelligence Laboratory*
*Department of Brain & Cognitive Sciences*
*Massachusetts Institute of Technology*
*Cambridge, MA 02139, USA*

**Pat Shafto**                                     P.SHAFTO@LOUISVILLE.EDU
**Baxter Eaves**                                   B0EAVE01@LOUISVILLE.EDU
*Department of Psychology*
*University of Louisville*
*Louisville, KY 40292, USA*

## Abstract

Is it possible to make statistical inference broadly accessible to non-statisticians without sacrificing mathematical rigor or inference quality? This paper describes BayesDB, a probabilistic programming platform that aims to enable users to query the probable implications of their data as directly as SQL databases enable them to query the data itself. This paper focuses on four aspects of BayesDB: (i) BQL, an SQL-like query language for Bayesian data analysis, that answers queries by averaging over an implicit space of probabilistic models; (ii) techniques for implementing BQL using a broad class of multivariate probabilistic models; (iii) a semi-parametric Bayesian model-builder that auomatically builds ensembles of factorial mixture models to serve as baselines; and (iv) MML, a "meta-modeling" language for imposing qualitative constraints on the model-builder and combining baseline models with custom algorithmic and statistical models that can be implemented in external software. BayesDB is illustrated using three applications: cleaning and exploring a public database of Earth satellites; assessing the evidence for temporal dependence between macroeconomic indicators; and analyzing a salary survey.

**Keywords:** Probabilistic programming, Bayesian inference, probabilistic databases, multivariate statistics, nonparametric Bayes, automatic machine learning

## 1. Introduction

Is it possible to make statistical inference broadly accessible to non-statisticians without sacrificing mathematical rigor or inference quality? This paper describes BayesDB, a system that enables users to query the probable implications of their data as directly as SQL

databases enable them to query the data itself. By combining ordinary SQL with three new primitives — SIMULATE, INFER, and ESTIMATE — users of BayesDB can detect predictive relationships between variables, retrieve statistically similar data items, identify anomalous data points and variables, infer missing values, and synthesize hypothetical subpopulations. The default modeling assumptions that BayesDB makes are suitable for a broad class of problems (Mansinghka et al., 2015; Wasserman, 2011), but statisticians can customize these assumptions when necessary. BayesDB also enables domain experts that lack statistical expertise to perform qualitative model checking (Gelman et al., 1995) and encode simple forms of qualitative prior knowledge.

BayesDB consists of four components, integrated into a single probabilistic programming system:

1. The Bayesian Query Language (BQL), an SQL-like query language for Bayesian data analysis. BQL programs can solve a broad class of data analysis problems using statistically rigorous formulations of cleaning, exploration, confirmatory analysis, and predictive modeling. BQL defines these primitive operations for these workflows in terms of Bayesian model averaging over results from an implicit set of multivariate probabilistic models.

2. A mathematical interface that enables a broad class of multivariate probabilistic models, called generative population models, to be used to implement BQL. According to this interface, a data generating process defined over a fixed set of variables is represented by (i) an infinite array of random realizations of the process, including any observed data, and (ii) algorithms for simulating from arbitrary conditional distributions and calculating arbitrary conditional densities. This interface permits many statistical operations to be implemented once, independent of the specific models that will be used to apply these operations in the context of a particular data table.

3. The BayesDB Meta-modeling Language (MML), a minimal probabilistic programming language. MML includes constructs that enable statisticians to integrate custom statistical models — including arbitrary algorithmic models contained in external software — with the output of a broad class of Bayesian model building techniques. MML also includes constructs for specifying qualitative dependence and independence constraints.

4. A hierarchical, semi-parametric Bayesian "meta-model" that automatically builds ensembles of generalized mixture models from database tables. These ensembles serve as baseline data generators that BQL can use for data cleaning, initial exploration, and other routine applications.

This design insulates end users from most statistical considerations. Queries are posed in a qualitative probabilistic programming language for Bayesian data analysis that hides the details of probabilistic modeling and inference. Baseline models can be built automatically and customized by statisticians when necessary. All models can be critically assessed and qualitatively validated via predictive checks that compare synthetic rows (generated via BQL's SIMULATE operation) with rows from the original data. Instead of hypothesis testing, dependencies between variables are obtained via Bayesian model selection.

BayesDB is "Bayesian" in two ways:

1. In BQL, the objects of inference are rows, and the underlying probability model forms a "prior" probability distribution on the fields of these rows. This is then constrained by row-specific observations to create a posterior distribution of field values. Without this prior, it would be impossible to simulate rows or infer missing values from partial observations.

2. In MML, the default meta-model is Bayesian in that it assigns a prior probability to a very broad class of probabilistic models and narrows down on probable models via Bayesian inference. This prior is unusual in that it encodes a state of ignorance rather than a strong inductive constraint. MML also provides instructions for augmenting this prior to incorporate qualitative and quantitative domain knowledge.

In practice, it is useful to use BQL for Bayesian queries against models built using non-Bayesian or only partially Bayesian techniques. For example, MML supports composing the default meta-model with modeling techniques specified in external code that need not be Bayesian. However, the default is to be Bayesian for both model building and query interpretation, as this ensures the broadest applicability of the results.

This paper focuses on the technical details of BQL, the data generator interface, the meta-model, and the MML. It also illustrates the capabilities of BayesDB using three applications: cleaning and exploring a public database of Earth satellites, discovering relationships in measurements of macroeconomic development of countries, and analyzing salary survey data. Empirical results are based on a prototype implementation that embeds BQL into `sqlite3`, a lightweight, open-source, in-memory database.

## 1.1 A conceptual illustration

This section illustrates data analysis using the MML and BQL on a synthetic example based on analysis of electronic health records. SQL databases make it easy to load data from disk and run queries that filter and retrieve the contents. The first step in using BayesDB is to load data that describes a statistical (sub)population into a table, with one row per member of the population, and one column per variable:

```
CREATE POPULATION patients WITH DATA FROM patients.csv;

SELECT age, has_heart_disease FROM patients WHERE age > 30 LIMIT 3;
```

| age | has_heart_disease |
|-----|-------------------|
| 66  | ???               |
| 44  | yes               |
| 31  | ???               |

Once data has been loaded, a *population schema* needs to be specified. This schema specifies the statistical characteristics of each example. For example, whether it is categorical or numerical, and if it is categorical, how many outcomes are there and how is each outcome represented. After an initial schema has been specified — using a mix of automatic inference

and manual specification — the schema can be customized using instructions in the Meta-modeling Language (MML).

```
GUESS POPULATION SCHEMA FOR patients;


ALTER POPULATION SCHEMA FOR patients
SET DATATYPE FOR num_hosp_visits TO COUNT;


CREATE DEFAULT METAMODEL FOR patients;
ALTER METAMODEL FOR patients ENSURE will_readmit DEPENDENT ON dialysis;


ALTER METAMODEL FOR patients
MODEL infarction GIVEN gender, age, weight, height, cholesterol, bp
USING CUSTOM MODEL FROM infarction_regression.py;
```

One distinctive feature of MML is that it includes instructions for *qualitative probabilistic programming*. These instructions control the behavior of the automatic modeling machinery in the MML runtime. In this example, these constraints include the assertion of a dependence between the presence of a chronic kidney condition and future hospital readmissions. They also include the specification of a custom statistical model for the `infarction` variable, illustrating one way that discriminative and non-probabilistic approaches to inference can be integrated into BayesDB.

The next step is to use the MML to build an ensemble of general-purpose models for the data, subject to the specified constraints:

```
INITIALIZE 100 MODELS FOR patients;
ANALYZE patients FOR 3 HOURS CHECKPOINT EVERY 10 MINUTES;
```

Each of these 100 models is a *generative population model* (GPM) that represents the joint distribution on all possible measurements of an infinite population with the given population schema. These models are initially drawn accordingt to a broad prior probability distribution over a large hypothesis space of possible GPMs. Until the observed data has been analyzed, BQL will thus report broad uncertainty for all its query responses.

Once the models are sufficiently adapted to the data, it is possible to query its probable implications. The following query quantifies over columns, rather than rows, and retrieves the probability of a marginal dependence between three (arbitrary) variables and `height`:

```
ESTIMATE COLUMN NAME, PROBABILITY OF DEPENDENCE WITH height
FROM COLUMNS OF patients LIMIT 3;
```

| column name | p( dep. with height ) |
|---|---|
| height | 1.0 |
| infarction | 0.08 |
| gender | 0.99 |

Point predictions can be accessed by using the `INFER` instruction, a natural generalization of `SELECT` from SQL:

```
INFER age, has_heart_disease FROM patients
WHERE age > 30 WITH CONFIDENCE 0.8 LIMIT 3;
```

| age | has_heart_disease |
|-----|-------------------|
| 66  | yes               |
| 44  | yes               |
| 31  | ???               |

In this example, only one of the missing values could be inferred with the specified confidence level. The probabilistic semantics of CONFIDENCE will be discussed later in this paper.

BQL also makes it straightforward to generate synthetic sub-populations subject to a broad class of constraints:

```
SIMULATE height, weight, blood_pressure FROM patients 3 TIMES
GIVEN gender = male AND age < 10
```

| height | weight | blood_pressure |
|--------|--------|----------------|
| 46     | 80     | 110            |
| 38     | 60     | 80             |
| 39     | 119    | 120            |

The SIMULATE operator gives BQL users access to samples from the posterior predictive distribution induced by the implicit underlying set of models. This is directly useful for predictive modeling and also decision-theoretic choice implemented using Monte Carlo estimation of expected utility (Russell and Norvig, 2003). It also enables predictive checking: samples from SIMULATE can be compared to the results returned by SELECT. Finally, domain experts can use SIMULATE to scrutinize the implications of the underlying model ensemble, both quantitatively and qualitatively.

## 2. Example Analyses

This section describes three applications of the current BayesDB prototype:

1. Exploring and cleaning a public database of Earth satellites.

2. Assessing the evidence for dependencies between indicators of global poverty

3. Analyzing data from a salary survey.

BQL and MML constructs are introduced via real-world uses; a discussion of their formal interpretation is provided in later sections.

### 2.1 Exploring and cleaning a public database of Earth satellites

The Union of Concerned Scientists maintains a database of 1000 Earth satellites. For the majority of satellites, it includes kinematic, material, electrical, political, functional, and economic characteristics, such as dry mass, launch date, orbit type, country of operator, and purpose. Here we show a sequence of interactions with a snapshot of this database using the bayeslite implementation of BayesDB.

### 2.1.1 Inspecting the data.

We start by loading the data and looking at a sample. This process uses a combination of ordinary SQL and convenience functions built into `bayeslite`. The first step is to create a population from the raw data:

```
CREATE POPULATION satellites FROM ucs_database.csv
```

One natural query is to find the International Space Station, a well-known satellite:

```
SELECT * FROM satellites WHERE Name LIKE 'International Space Station%'
```

| Variable | Value |
|----------|-------|
| Name | International Space Station (ISS ...) |
| Country_of_Operator | Multinational |
| Operator_Owner | NASA/Multinational |
| Users | Government |
| Purpose | Scientific Research |
| Class_of_Orbit | LEO |
| Type_of_Orbit | Intermediate |
| Perigee_km | 401 |
| Apogee_km | 422 |
| Eccentricity | 0.00155 |
| Period_minutes | 92.8 |
| Launch_Mass_kg | NaN |
| Dry_Mass_kg | NaN |
| Power_watts | NaN |
| Date_of_Launch | 36119 |
| Anticipated_Lifetime | 30 |
| Contractor | Boeing Satellite Systems (prime)/Multinational |
| Country_of_Contractor | Multinational |
| Launch_Site | Baikonur Cosmodrome |
| Launch_Vehicle | Proton |
| Source_Used_for_Orbital_Data | www.satellitedebris.net 12/12 |
| longitude_radians_of_geo | NaN |
| Inclination_radians | 0.9005899 |

This result row illustrates typical characteristics of real-world databases such as heterogeneous data types and missing values.

### 2.1.2 Building baseline models.

Before exploring the implications of the data, it is necessary to obtain a collection of probabilistic models. The next two MML instructions produce a collection of 16 models, using roughly 4 minutes of analysis total.

```
INITIALIZE 16 MODELS FOR satellites;
ANALYZE satellites FOR 4 MINUTES WAIT;
```

Each of the 16 models is a separate GPM produced by an independent Markov chain for approximate posterior sampling in the default semi-parametric factorial mixture meta-model described earlier. This number of models and amount of computation is typical for the exploratory analyses done with our prototype implementation; this is sufficient for roughly 100 full sweeps of all latent variables.

### 2.1.3 ANSWERING HYPOTHETICALS.

The satellites database should in principle inform the answers to a broad class of hypothetical or "what if?" questions. For example, consider the following question:

> *Suppose you receive a report indicating the presence of a previously undetected satellite in geosynchronous orbit with a dry mass of 500 kilograms. What countries are most likely to have launched it, and what are its likely purposes?*

Answering this question requires knowledge of satellite engineering, orbital mechanics, and the geopolitics of the satellite industry. It is straightforward to answer this question using BQL. The key step is to generate a synthetic population of satellites that reflect the given constraints:

```
SIMULATE country_of_operator, purpose FROM satellites
GIVEN Class_of_orbit = GEO, Dry_mass_kg = 500 LIMIT 1000;
```

Figure 1 shows the results of a simple aggregation of these results, counting the marginal frequencies of various countries and purposes and sorting accordingly. The most probable explanation, carrying roughly 25% of the probability mass, is that it is a communications satellite launched by the USA. It is also plausible that it might have been launched other major space powers such as Russia or China, and that it might have a military purpose.

The satellites data are too sparse and ambiguous for frequency counting to be a viable alternative. Consider an approach based on finding satellites that match the discrete GEO constraint and are within some ad-hoc tolerance around the observed dry mass:

```
SELECT country_of_operator, purpose, Class_of_orbit, Dry_mass_kg
FROM satellites
WHERE Class_of_orbit = "GEO" AND Dry_Mass_kg BETWEEN 400 AND 600;
```

This SQL query returns just 2 satellites, both Indian:

|   | Country _ of _ Operator | Purpose | Class _ of _ Orbit | Dry _ Mass _ kg |
|---|---|---|---|---|
| 0 | India | Communications | GEO | 559 |
| 1 | India | Meteorology | GEO | 500 |

Presuming our intuition about satellite mass is flawed, we might issue another query to look at a broader range of satellites:

```
SELECT country_of_operator, purpose, Class_of_orbit, Dry_mass_kg
FROM satellites
WHERE Class_of_orbit = 'GEO' AND Dry_Mass_kg BETWEEN 300 AND 700
```
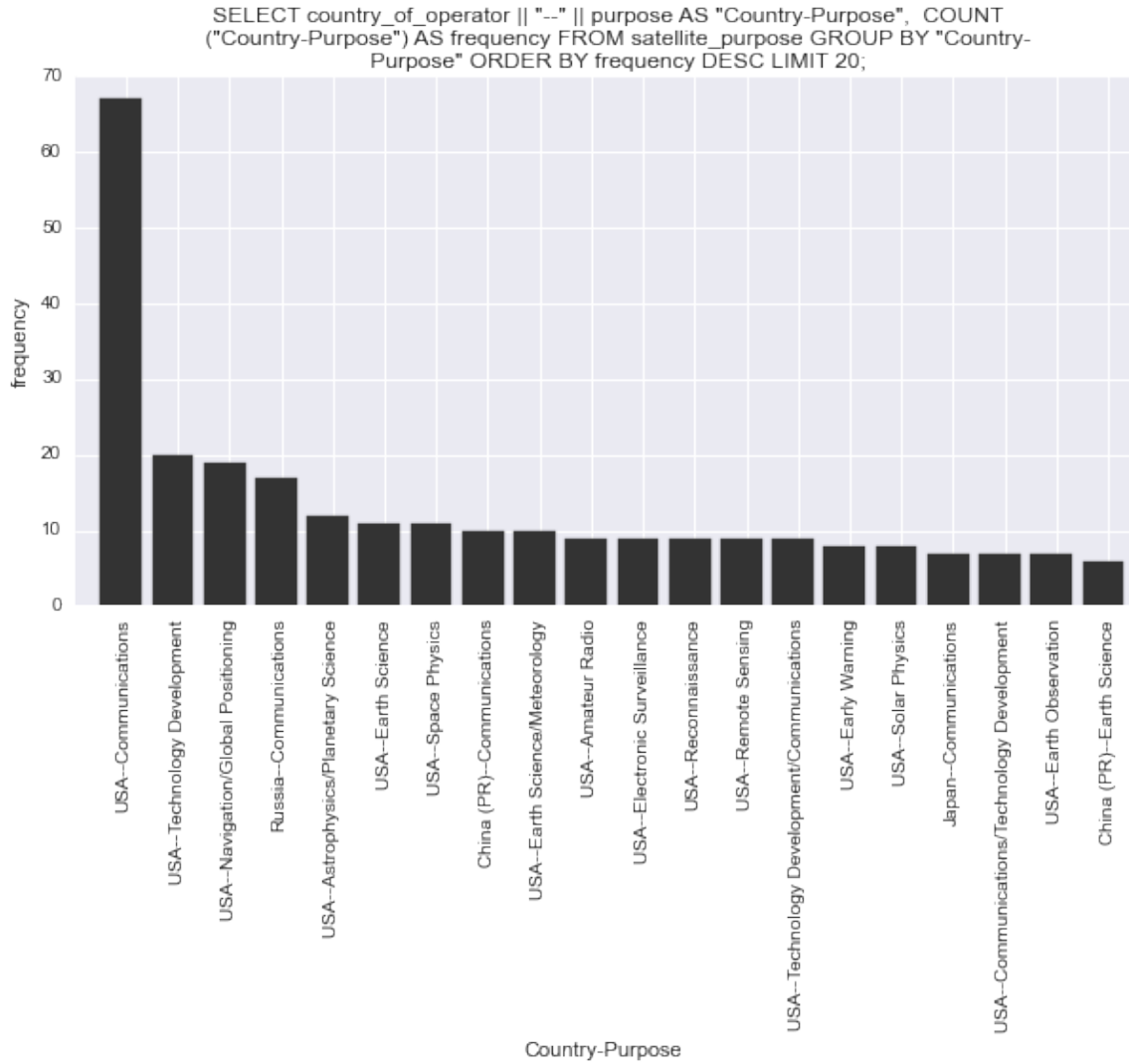
Figure 1: **The most probable countries and purposes of a satellite with a 500 kilogram dry mass in geosynchronous orbit.** See main text for discussion.

The results still do not give any real insight into the likely purpose of this satellite:

|    | Country_of_Operator | Purpose | Class_of_Orbit | Dry_Mass_kg |
|----|----|----|----|----|
| 0  | Malaysia | Communications | GEO | 650 |
| 1  | Israel | Communications | GEO | 646 |
| 2  | Luxembourg | Communications | GEO | 700 |
| 3  | Russia | Communications | GEO | 620 |
| 4  | China (PR) | Earth Science | GEO | 620 |
| 5  | China (PR) | Earth Science | GEO | 620 |
| 6  | China (PR) | Earth Science | GEO | 620 |
| 7  | India | Communications | GEO | 559 |
| 8  | India | Navigation | GEO | 614 |
| 9  | India | Meteorology | GEO | 500 |
| 10 | Malaysia | Communications | GEO | 650 |
| 11 | Multinational | Earth Science/Meteorology | GEO | 320 |
| 12 | United Kingdom | Communications | GEO | 660 |
| 13 | Norway | Communications | GEO | 646 |

Without deep expertise in satellites, and significant expertise in statistics, it is difficult to know whether or not these results can be trusted. How does the set of satellites vary as the thresholds on `Dry_Mass_kg` are adjusted? How locally representative and comprehensive is the coverage afforded by the data? Are there indirect, multivariate dependencies that ought to be taken into account, to determine which satellites are most similar? How should existing satellites be weighted to make an appropriate weighted sample against which to calculate frequencies? In fact, small modifications to the tolerance on `Dry_Mass_kg` yield large changes in the result set.

### 2.1.4 IDENTIFYING PREDICTIVE RELATIONSHIPS BETWEEN VARIABLES.

A key exploratory task is to identify those variables in the database that appear to predict one another. This is closely related to the key confirmatory analysis question of assessing the evidence for a predictive relationship between any two particular variables.

To quantify the evidence for (or against) a predictive relationship between two pairs of variables, BQL relies on information theory. The notion of dependence between two variables A and B is taken to be mutual information; the amount of evidence for dependence is then the probability that the mutual information between A and B is nonzero. If the population models are obtained by posterior inference in a meta-model — as is the case with MML — then this probability approximates the posterior probability (or strength of evidence) that the mutual information is nonzero.

```
ESTIMATE DEPENDENCE PROBABILITY FROM PAIRWISE COLUMNS OF satellites;
```

Figure 2 shows the results from this query. There are several groups of variables with high probability of mutual interdependence. For example, we see a definite block of geopolitically related variables, such as the country of contractor & operator, the contractor's identity,

and the location of the satellite (if it is in geosynchronous orbit). The kinematic variables describing orbits, such as perigee, apogee, period, and orbit class, are also shown as strongly interdependent. A domain expert with sufficiently confident domain knowledge can use this overview of the predictive relationships to assess the value of the data and the validity of the baseline models.

It is also instructive to compare the heatmap of pairwise dependence probabilities with alternatives from statistics. Figure 2 also shows heatmap that results from datatype-appropriate measures of correlation. The results from correlation are sufficiently noisy that it would be difficult to trust inferences from techniques that use correlation to select variables. Furthermore, the most causally unambiguous relationships, such as the kinematic constraints relating perigee, apogee, and orbital period, not detected by correlation.

### 2.1.5 DETECTING MULTIVARIATE ANOMALIES.

Another key aspect of exploratory analysis is identifying anomalous values, including both (univariate) outliers and multivariate anomalies. Anomalies can arise due to errors in data acquisition, bugs in upstream preprocessing software (including binning of continuous variables or translating between different discrete outcomes), and runtime failures. Anomalies can also arise due to genuine surprises or changes in the external environment.
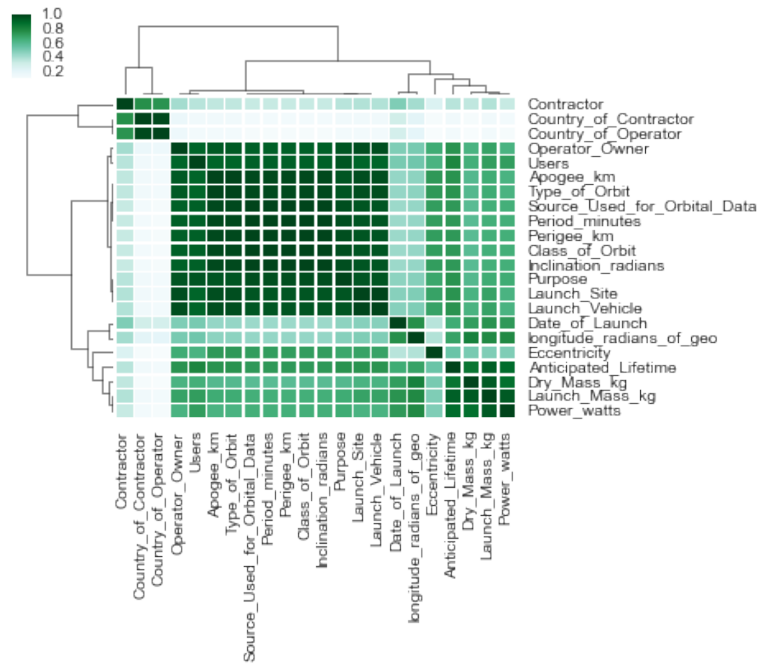
Using BQL, multivariate anomalies can be detected by assessing the predictive probability density of each measurement, and ordering from least to most probable. Here we illustrate this using a simple example: ordering the geosynchronous satellites according to the probability of their recorded orbital period:

```
ESTIMATE name, class_of_orbit, period_minutes AS TAU,
PREDICTIVE PROBABILITY OF period_minutes AS "Pr[TAU]"
FROM satellites
ORDER BY ``Pr[TAU]'' ASCENDING LIMIT 10
```
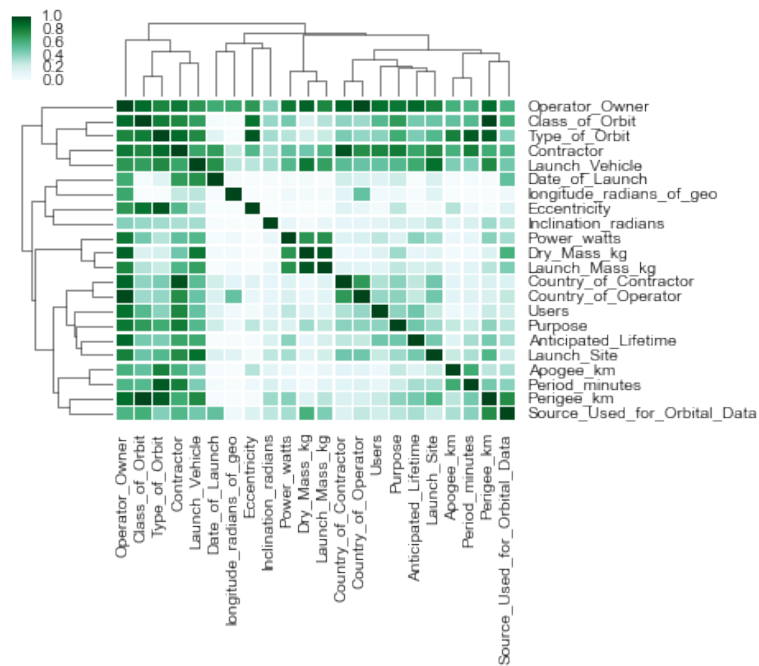
This BQL query produces the following table of results:

|   | Name | Class_of_Orbit | TAU | Pr[TAU] |
|---|---|---|---|---|
| 0 | AEHF-3 (Advanced Extremely High Frequency sate... | GEO | 1306.29 | 0.001279 |
| 1 | AEHF-2 (Advanced Extremely High Frequency sate... | GEO | 1306.29 | 0.001292 |
| 2 | DSP 20 (USA 149) (Defense Support Program) | GEO | 142.08 | 0.002656 |
| 3 | Intelsat 903 | GEO | 1436.16 | 0.003239 |
| 4 | BSAT-3B | GEO | 1365.61 | 0.003440 |
| 5 | Intelsat 902 | GEO | 1436.10 | 0.003492 |
| 6 | SDS III-6 (Satellite Data System) NRO L-27, Gr... | GEO | 14.36 | 0.003811 |
| 7 | Advanced Orion 6 (NRO L-15, USA 237) | GEO | 23.94 | 0.003938 |
| 8 | SDS III-7 (Satellite Data System) NRO L-38, Dr... | GEO | 23.94 | 0.003938 |
| 9 | QZS-1 (Quazi-Zenith Satellite System, Michibiki) | GEO | 1436.00 | 0.004446 |

Recall that a geosynchronous orbit should take 24 hours or 1440 minutes. Rows 7 and 8 appear to be unit conversion errors (hours rather than minutes). Rows 2 and 6 appear to be decimal placement errors. Note that row 2 is not an outlier: some satellites have an orbital

(a) A heatmap depicting the pairwise probabilities of dependence between all pairs of variables. The rows and columns are both permuted according to a single ordering obtained via agglomerative hierarchical clustering to highlight multivariate interactions.



(b) The pairwise correlation matrix; note that many causal relationships are not detected by simple correlations. See main text for discussion.

Figure 2: **Detecting predictive relationships in the satellites database.**

period of roughly two hours. It is only anomalous in the context of the other variables that are probably predictive of orbital period, such as orbit class.

2.1.6 Inferring missing values.

A key application of predictive modeling is inferring point predictions for missing measurements. This can be necessary for cleaning data before downstream processing. It can also be of intrinsic interest, e.g. in classification problems. The satellites database has many missing values. Here we show an `INFER` query that infers missing orbit types and returns both a point estimate and the confidence in that point estimate:

```
INFER EXPLICIT
anticipated_lifetime, perigee_km, period_minutes, class_of_orbit,
PREDICT type_of_orbit AS inferred_orbit_type CONFIDENCE inferred_orbit_type_conf
FROM satellites
WHERE type_of_orbit IS NULL;
```

This form of `INFER` uses the `EXPLICIT` modifier that exposes both predicted values and their associated confidence levels to be included in the output. Figure 3 shows a visualization of the results. The panel on the bottom left shows that the confidence depends on the orbit class and on the predicted value for the inferred orbit type. For example, there is typically moderate to high confidence for the orbit type of `LEO` satellites — and high confidence (but some variability in confidence) for those with `Sun-Synchronous` orbits. Satellites with `Elliptical` orbits may be assigned a `Sun-Synchronous` type with moderate confidence, but for other target labels confidence is generally lower. After examining the overall distribution on confidences, it can be natural to filter `INFER` results based on a manually specified confidence threshold. Note that many standard techniques for imputation from statistics correspond to `INFER ... WITH CONFIDENCE 0`.

2.1.7 Integrating a kinematic model for elliptical orbits.

Can we improve over the baseline models by integrating causal knowledge about satellites? MML can be used to compose GPMs built by the default model builder with algorithmic and/or statistical models specified as external software. Here we integrate a simple model for elliptical orbits:

```
ALTER METAMODEL FOR satellites
MODEL perigee_km, apogee_km GIVEN period_minutes, eccentricity
USING CUSTOM MODEL FROM stochastic_kepler.py;
ANALYZE FOREIGN PREDICTORS FOR 1 MINUTE;
```

The underlying foreign predictor implements Kepler's laws:

$$R_{min} = \tau^{\frac{2}{3}}(1.0 - \epsilon) - R_{GEO}$$
$$R_{max} = \tau^{\frac{2}{3}}(1.0 + \epsilon) - R_{GEO}$$
$$X_{(r^*,\texttt{apogee\_km})} \sim N(R_{min}, \sigma^2)$$
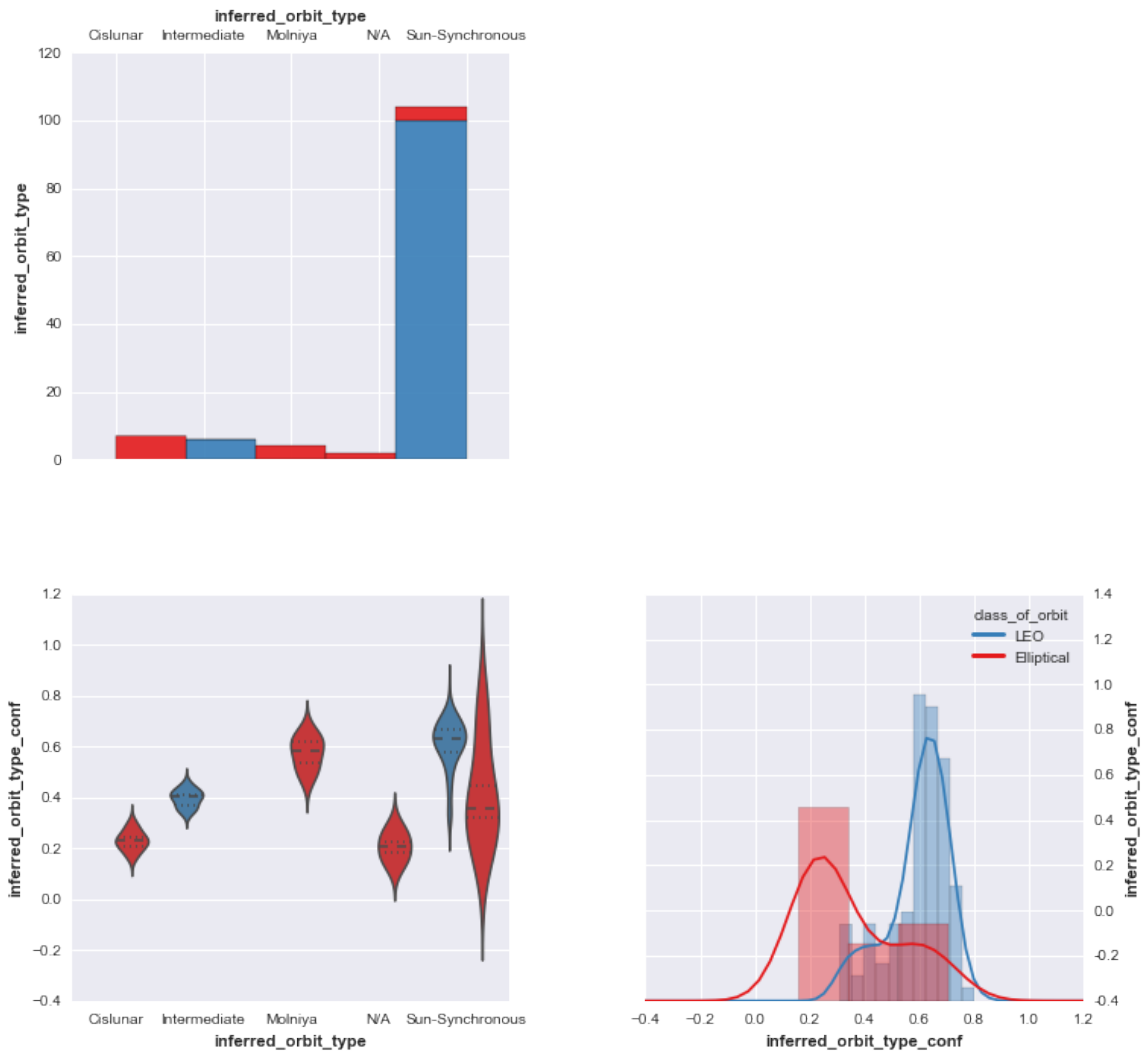$$X_{(r^*,\texttt{perigee\_km})} \sim N(R_{max}, \sigma^2)$$

Figure 3: **A visualization of inferred point estimates for `type_of_orbit` and the confidence in those point estimates.** See main text for discussion.

Here, $\epsilon$ is set via the `eccentricity` measurement for row $r^*$, and $\tau$ is set via the `period_minutes` measurement. $R_{GEO}$ is a fixed constant inside the foreign predictor representing the radius of the Earth in kilometers. $\sigma$ is a parameter that determines the noise and is set when the variable subset for the foreign predictor instantiation is `ANALYZE`d. Note that foreign predictors are essentially GPMs and accordingly must implement generic simulate(...) and logpdf(...) methods. For algorithmic forward models with numerical outputs, MML provides a default wrapper that uses importance sampling with resampling to approximately generate conditional samples and estimate marginal densities. This is how Kepler's laws — in the form of a forward simulation — are turned into a generative model for kinematic variables that can be conditionally simulated in arbitrary directions.

Figure 4 and Figure 5 show the results. A detailed discussion of the relative merits of empirical versus analytical modeling is beyond the scope of this paper. However, it is clear that neither the empirical approach nor the analytical approach is universally dominant. The empirical approach is able to correctly locate the empirical probability mass — including multiple modes — but underfits. The orbital mechanics approach yields inferences that are typically in much tighter accord with the kinematic data. This is unsurprising: these are the patterns of covariation that led to the development of quantitative models and "hard" natural sciences. However, there are many satellites for which Kepler's laws are not in accord with the data. This reflects many factors, including data quality errors as well as legitimate gaps between idealized mathematical laws and fine-grained empirical records of real-world phenomena. For example, at the time Kepler's laws were formulated, orbiting bodies lacked engines.

### 2.1.8 Combing random forests, causal models, and nonparametric Bayes.

Because MML supports model composition, it is straightforward to build hybrid models that integrate techniques from subfields of machine learning that might seem to be in conflict. Figure 6 shows the transcript of a complete MML session that builds such a hybrid model. Random forests are used to classify orbits into types; Kepler's Laws are used to relate period, perigee, and apogee; and the default semi-parametric Bayesian meta-model is used for all remaining variables (with two variables coming with overridden datatypes).

Later sections of this paper explain how these three modeling approaches are combined to answer individual BQL queries.

## 2.2 Assessing the evidence for dependencies between indicators of global poverty

In the early 21st century it is widely believed that resolving extreme poverty around the world will be accomplished by empowering individuals to resolve their poverty. Governments and NGOs encourage this process through a variety of interventions, many of them combining material assistance with policy changes. In principle, policies should be driven by quantitative data-driven understanding of international economic development. In practice, international economic data is sparse, unreliable, and highly aggregated. These data limitations create substantial obstacles to understanding the situational context of successful or unsuccessful interventions and policies.

```
SIMULATE period_minutes, apogee_km FROM satellites_kepler LIMIT 100;
```
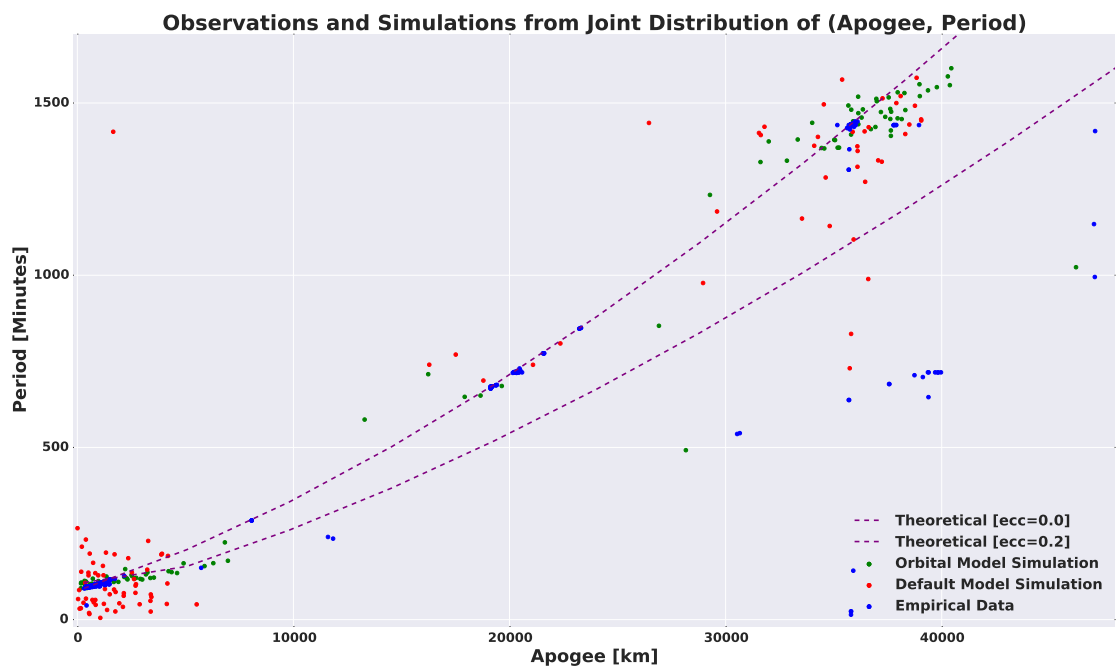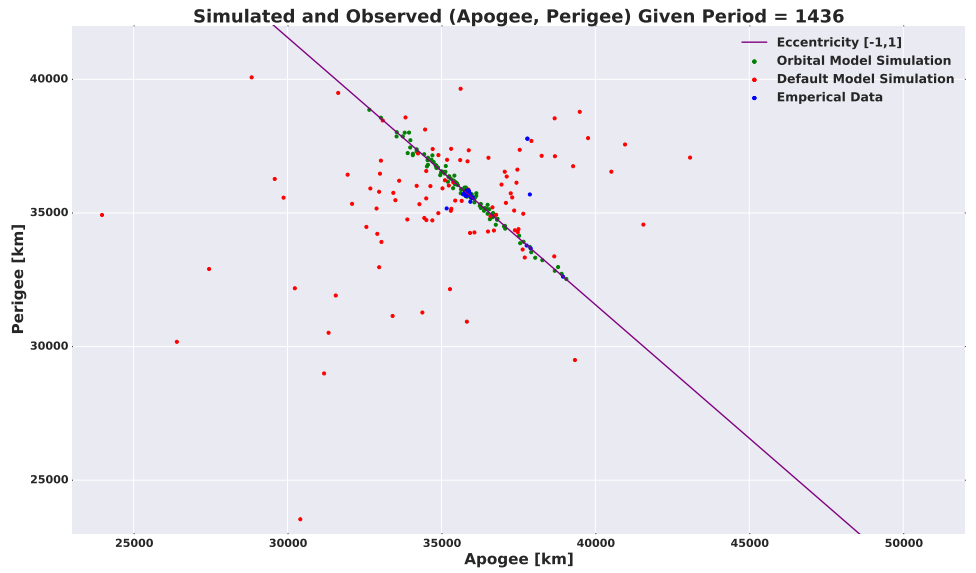


Figure 4: **Integrating empirical baseline models with a noisy version of Kepler's laws.** Neither the empirical approach nor the analytical approach is universally dominant in terms of accuracy. See main text for discussion.

SIMULATE perigee_km, apogee_km FROM satellites_kepler ASSUMING
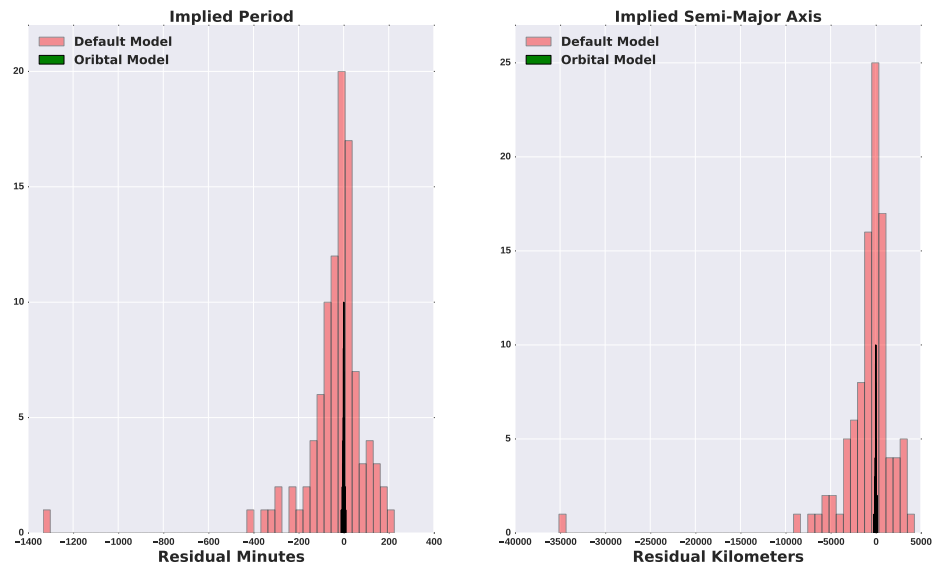period_minutes = 1436 LIMIT 100;

Figure 5: **Comparing conditional predictions of an empirical model with Kepler's laws.** See main text for discussion.

```
CREATE POPULATION satellites
  FROM ucs_satellites.csv

CREATE METAMODEL sat_keplers ON satellites
  USING composer(
  random_forest (
    Type_of_Orbit (CATEGORICAL)
      GIVEN Apogee_km, Perigee_km,
            Eccentricity, Period_minutes,
            Launch_Mass_kg, Power_watts,
            Anticipated_Lifetime,
 Class_of_orbit
   ),
   foreign_model (
     source = 'keplers_laws.py',
     Period_Minutes (NUMERICAL)
       GIVEN Perigee_km, Apogee_km
   ),
   default_metamodel (
     Country_of_Operator CATEGORICAL,
     Inclination_radians NUMERICAL
   )
);

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

Figure 6: A complete MML session that builds a hybrid model integrating techniques from subfields of machine learning that might seem to be in conflict.

The "Gapminder" data set, collected and curated by Hans Rosling at the Karolinska Institutet, is the most well known and extensive sets of longitudinal global developmental indicators. Representing over 500 indicators, 400 countries, and 500 years of data, it covers the colonial era, industrial revolution, socio-political upheavals around the world in the 20th century, and the first decade of the 21st. Containing over 2 million observations, the data has been used as the basis for a compelling set of data animations and the most widely viewed TED talk on statistics.

To date, analysis of this data has been minimal, as it requires intensive preprocessing and cleaning. Different analytical methdologies require different approaches to imputation and variable selection; as a result, results from different teams are difficult to compare. Here we show how to explore the data with BayesDB and assess the evidence for predictive relationships between different macroeconomic measures of development.

### 2.2.1 Exploring the Data with SQL

The raw form of the data is ∼500 Excel spreadsheets, each containing longitudinal data for ∼300 countries over ∼100 years. However, the dataset only contains ∼2 million observations, i.e. 97% of the data is missing. Figure 7 shows key indicators of the data around size, missing records, and the relationship between data availability and countries, records, and years. The primary data is mmodeled in SQL as a âĂIJfactâĂİ table structure. This relatively-normalized representation easily models the sparse matrix and allows us to use a combination of SQL and Python data science tools to craft our population structures.

The histograms in Figure 7 show the breadth and also the variability of the data. The histogram by year in Figure 7a shows that data is complete for only recent history, and in fact that some predicted data continues into the future, and that data for some indicators is only available every 10 years. The histogram by country in Figure 7a shows that data availability varies by country (the most described country is Sweden), that many countries have reasonably complete data, but that there is a long tail of countries with sparse data, including countries that no longer exist and with inconsistent or disputed naming. Figure 7c shows that there is also a variance by indicators, because different measurements are collected by different agencies with different expectations and data policies.

The data has already been subject to extensive visualization and descriptive analytics by the Gapminder project. This paper focuses on the use of MML to model the data and BQL to query its probable implications.
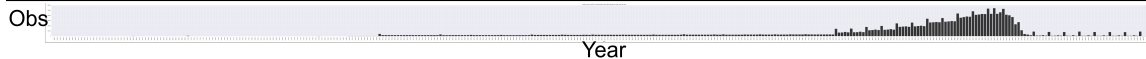
### 2.2.2 Detecting Basic and Longitudinal Dependence

Our analysis focuses on the 53 variables with most complete data for the years 1999-2008. It is straightforward to create an ensemble of models for this subset:

```
GUESS POPULATION SCHEMA FOR dense_gapminder;
INITIALIZE 64 MODELS FOR dense_gapminder;
ANALYZE todo FOR 300 MINUTES WAIT;
```

The probability of dependence heatmap that results is shown in Figure 8. Indicators such as total population and urban percentage form blocks containing their values for all 10 years contained in the dataset. This shows that the default GPM was able to extract the
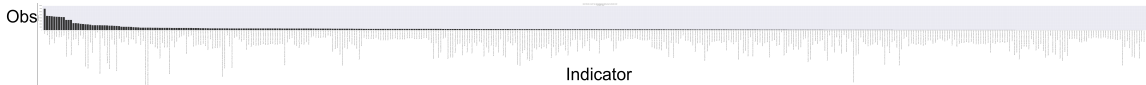
| observation_count | indicator_count | country_count | year_count | coverage |
|:---:|:---:|:---:|:---:|:---:|
| 2082193 | 464 | 405 | 364 | 0.03044 |



(a) **Observation volume by year, 1086-2100.** Note that the most complete data is for only recent history.



(b) **Observation volume by country.** Note that the data becomes very limited for some countries, and includes countries that no longer exist.



(c) **Observation volume by indicator.** Note that some indicators are much more complete or extensive in terms of year and country than others.

Figure 7: **Gapminder data volume measures.** The dataset contains longitudinal records of ~500 macroeconomic variables for ~300 countries, spanning a ~100 year period. However, roughly 97% of the data is missing.
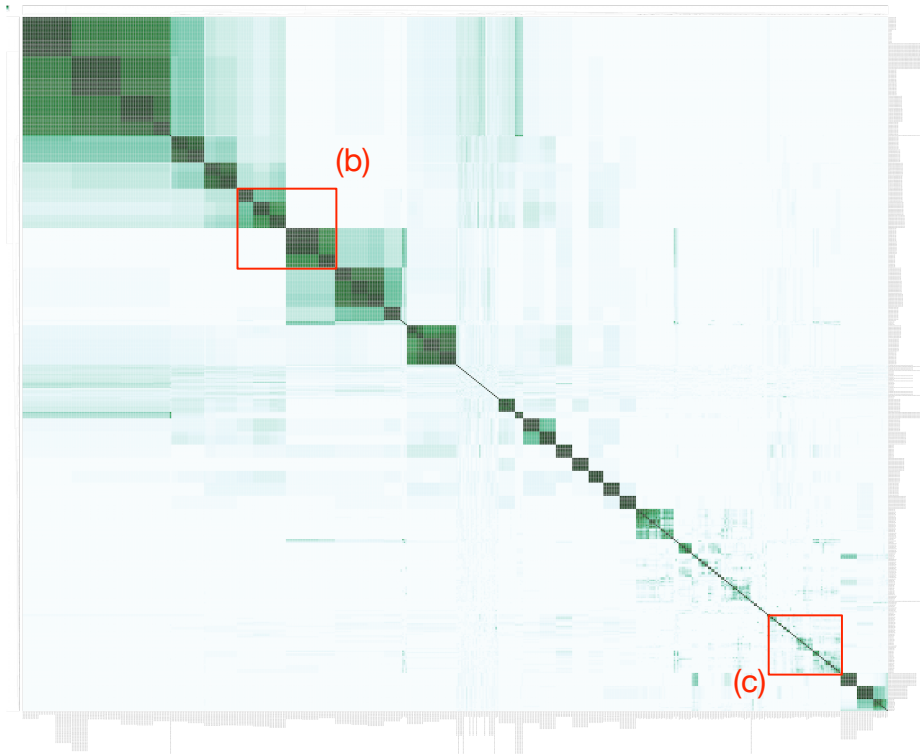
temporal dependence in these indicators. In other cases, such as measurements of the number of people killed in floods, the year to year dependence is much weaker. The heatmap also shows dependence between indicators, such as the block in the top right corner combining indicators of stress, urbanization, and fertility rate. Finally, it segregates data according to type sof indicators, as can be seen in Figure 8b where there is a sharp break from total measurements to per-capital.

If we analyze just the data for the year 2002, using 32 models for 3 minutes, we get a heatmap that highlights the dependence and independence between indicators. Figure 9 shows the details.
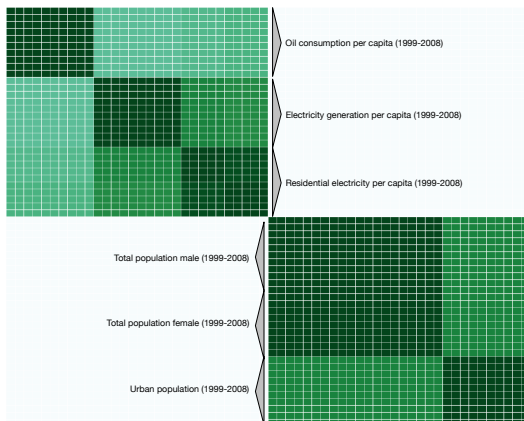
### 2.2.3 Measuring the Similarity of Countries

In order to help with the delivery of international aid and the design and analysis of interventions, decision makers often want a richer understanding of the similarities between countries. With BQL we can formulate these queries in general or against specific attributes. Figure 10 shows country similarities for different indicators. As expected, changing the indicator of interest can produce a very different similarity structure. Analyses that presume a single global similarity measure cannot pick up this context-specific structure.
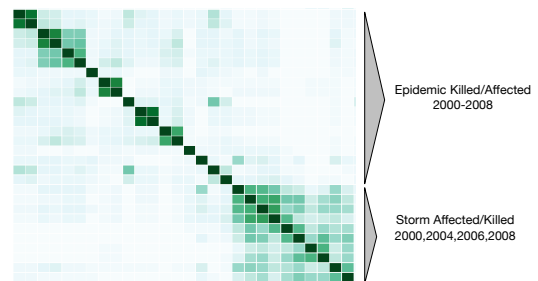
The authors are involved in an ongoing research partnership with the Bill and Melinda Gates Foundation aimed at integrating the Gapminder data with other relevant sources, including qualitative knowledge from domain eperts, and using it to drive empirically grounded policy and aid interventions.

(a) Probability of dependence heatmap for 40 indicators over 10 years.



(b) Indicators form 10 year runs, with a sharp break from totals to per capita.



(c) Natural disaster indicators cluster but do not have strong year-to-year dependence.

Figure 8: **Probability of dependence heatmap for the Gapminder data.** BayesDB detects temporal dependence within some indicators but not others, as well as dependence between some indicators (but not others). See main text for discussion.

Figure 9: **Probability of dependence for 40 indicators in 2002.** Of particular note are the blocks for population growth rates, natural disasters, HIV, total energy usage, and per capita energy usage.

(a) Countries enumerated in decending order of similarity to Kenya on Total Fertility Rate.



(b) Countries enumerated in descending order of similarity to the United States on Total Fertility Rate.



(c) Countries enumerated in decending order of similarity to Kenya on per capita energy production.



(d) Countries enumerated in descending order of similarity to the United States on per capita energy production.

Figure 10: **Examining the similarity of countries.** Kenya and the United States are similar to different countries, and the similarity structure with respect to per capita energy production and total fertility rate are significantly different.

(a) Probability of dependence between columns with default metamodel settings.

(b) Probability of dependence between columns when `state` is dependent on salary columns.

Figure 11: **Probability of dependence heat maps with and without dependence assertions in MML.** Note that in the second figure stronger dependencies were resolved overall.

## 2.3 Analyzing a salary survey

Surveys are a common source of multivariate data and a potentially appealing application for BayesDB. Here we show a preliminary analysis of a web-administered anonymous salary survey. Participants shared their compensation details along with information about their title, years of service, acheivements, employer, and geography.

## 2.4 Controlling Models with Qualitative Assumptions

This salary population provides an instructive example of applying qualitiative assumptions to a model. In this case, the first analysis of compensation data finds that geographic location (state, region) is not a factor in compensation. Domain experts suggest that is implausible, that cost of living and the competitive market in different cities is a significant factor in compensation of the survey participants. The following code can be used to apply this qualitative assumption:

```
ALTER METAMODEL FOR salary ENSURE total, equity, base, bonus DEPENDENT
ON state;
```

Without asserting the dependence, state is inferred to be dependent on region and independent of performance. After asserting a qualitative constraint, the probability of depencence heat map changes. Not only are the squares implied by that depencence colored to 1.0,

23

but other columns have re-aligned in their modeling. In particular, given this assumption, there appears to be more evidence of dependence between the 2012 and 2013 measures and core indicators such as years in the job, bonus, the presence of an equity stake, etc. Also, there appears to be less evidence that job title impacts the key compensation variables.

## 3. The Bayesian Query Language

The Bayesian Query Language (BQL) formalizes Bayesian data analysis without exposing the end user to model parameters, priors, and posteriors. This section describes the statistical operations that are implemented by the core BQL instruction set. It also describes the modeling formalism that is used to implement BQL.

### 3.1 Generative Population Models

BQL programs are executed against a weighted collection of *generative population models* (GPMs). At present, GPMs can be built in two ways:

1. Specified directly as external software libraries.

2. Inferred from data via probabilistic inference in a meta-model written in BayesDB's Meta-modeling Language.

GPMs can respond to queries about the joint distribution of the underlying data generating process as a whole or about the predictive distribution for a specific member of the population. The population can be thought of as a table, where individual members are specified by row indexes.

Each GPM induces a random table with a finite number of columns and an infinite number of rows, where each cell contains a random variable. BQL treats each BayesDB generator as a model of the data generating process underlying its associated table of observations. It is sometimes useful to query a GPM about hypothetical members of the population. This can be performed by using a row whose index $r^*$ may not be associated with any actual member; this can be guaranteed by generating a unique row index.

Each GPM is described by a *schema* $\mathcal{S}$ that must be compatible with the population schema for the population to which it is being applied. This schema is a tuple containing $(\texttt{typed-outputs}, \texttt{typed-inputs}, \texttt{body})$. The $\texttt{typed-outputs}$ component specifies the column indexes and statistical types of each column that the data generator will be responsible for producing. The $\texttt{typed-inputs}$ component specifies the indexes and statistical types of each column that the data generator can read from. The $\texttt{body}$ is an opaque binary that contains any GPM-specific configuration information, such as a probabilistic program.

Mathematically, the internals of a GPM $\mathcal{G} = (\Theta, \mathcal{Z}, O)$ consists of three parts:

1. Measurement-specific latent variables $\mathcal{Z} = \cup z_{(r,c)}$.

   There may be overlap between the latent variables for different measurements. If a GPM cannot track dependencies internally — or if it is based on a model class in which all measurements are coupled — then $z_{(r,c)} = \mathcal{Z}$.

2. Population-level latent variables $\Theta$.

   These are all latent variables that remain well-defined in the absence of all measurements. Examples include hyper-parameters and mixture component parameters.

3. Observations $O = \{(r_i, c_i, x_{(r_i, c_i)})\}$.

   These correspond to the observed measurements.

For example, a naive Bayesian GPM lacks any measurement-specific variables, i.e. $z_{(r,c)} = \emptyset$, and is completely characterized by a single vector of parameters $\Theta = \vec{\theta}_c$ for the probability models for each column. A finite mixture GPM would have $z_{(r,c)} = \{z_r\}$ be the cluster assignment for each row, and have $\Theta = \{\theta_{(c,l)} | l \in \mathcal{Z}\}$ be the component model parameters for each cluster.

Generative population models are required to satisfy the following conditional independence constraint:

$$x_{(r,c)} | \Theta, z_{(r,c)} \perp\!\!\!\perp x_{(r',c')} | \Theta, z_{(r',c')} \text{ unless } (r,c) = (r',c')$$

Note in particular that the observations $O$ need not be conditioned on directly, given $\Theta$ and $z_{(r,c)}$. This formalizes the requirement that the dependencies between the measurements in the population are completely mediated by the population-level latent variables and all relevant measurement-specific latent variables. In general, no other independence constraints are enforced by the interface. GPMs can thus be built around dense, highly-coupled model families such as low-dimensional latent spaces and convolutional neural networks.

### 3.1.1 AN INTERFACE TO GENERATIVE POPULATION MODELS

A GPM must implement the following interface:

1. $\mathcal{G} = \{\Theta, \mathcal{Z}\} = \text{initialize}( \text{schema} = \mathcal{S} )$

   Initialize a data generator with the given schema and return the resulting data generator $\mathcal{G}$. It ensures that storage has been allocated for the random variables $\Theta$ and $\mathcal{Z}$, storing the global latent variables and the local latent variables, respectively.

2. $\vec{s}_i = \text{simulate}(\mathcal{G}, \text{givens} = \{(r_j, c_j, x_{(r,c_j)})\}, \text{targets} = \{r_k, c_k\}, N)$

   Generate $N$ sampled values $\{\vec{s}_i\}$ from the specified distribution:

   $$\{\vec{s}_i\} \sim \{X_{(r_k, c_k)}\} | \{X_{(r_j, c_j)} = x_{(r_j, c_j)}\}, \Theta, \{z_{r,c} | (r,c) \in \{r_j, c_j\} \cup \{r_k, c_k\}\}$$

   The set of valid distributions includes all finite-dimensional joint distributions obtainable by conditioning on the values of arbitrary measurements and marginalizing over another arbitrary set.

3. $\log p = \text{logpdf}(\mathcal{G}, \text{givens} = \{(r_j, c_j, x_{(r_j, c_j)})\}, \text{query} = \{(r_k, c_k, q_{(r, c_k)})\})$

   Evaluate the log probability density of the specified conditional/marginal distribution at a target point:

   $$\log p = \log p(\{X_{(r_k, c_k)} = q_{(r_k, c_k)}\} | \{X_{(r_j, c_j)} = x_{(r_j, c_j)}\}, \Theta, \{z_{r,c} | (r,c) \in \{r_j, c_j\} \cup \{r_k, c_k\}\})$$

4. $d = \text{kl-divergence-given-G}(\mathcal{G}, \text{measurements\_A} = \{(r_i, c_i^a)\}, \text{measurements\_B} = \{(r_j, c_j^b)\},$
$\text{conditions\_C} = \{(r_k, c_k^c, x^k)\})$

This estimates the KL divergence of the set of measurements $A$ from the set of measurements $B$, conditioned on the given constraints $C$. KL calculations are central to model-independent data analysis. For example, to detect predictive relationships, it suffices to check for non-zero mutual information, which can be reduced to calculating the KL between the joint distribution over two variables and the product of the marginals.

It is included in the GPM interface because that allows a GPM implementer to supply an optimized implementation. Where such an implementation is not available, the KL can be estimated via simple Monte Carlo estimation:

$$D_{KL}^{\mathcal{G}}(\{X_{(r_i, c_i^a)}\}, \{X_{(r_j, c_j^b)}\}) = \sum_{\{x_i\} \in dom(\{X_{(r_i, c_i^a)}\})} p\left(\{X_{(r_i, c_i^a)}\} = \{x_i\} | \mathcal{G}\right) log\left(\frac{p\left(\{X_{(r_j, c_j^b)}\} = \{x_i\} | \mathcal{G}\right)}{p\left(\{X_{(r_i, c_i^a)}\} = \{x_i\} | \mathcal{G}\right)}\right)$$

$$\approx \sum_{\{x_i\}^k} log\left(\frac{p\left(\{X_{(r_j, c_j^b)}\} = \{x_i\}^k | \mathcal{G}\right)}{p\left(\{X_{(r_i, c_i^a)}\} = \{x_i\}^k) | \mathcal{G}\right)}\right)$$

$$\text{with} \{x_i\}^k \sim \{X_{(r_i, c_i^a)}\}$$

This interface is intentionally quite general. It needs to support an open set of primitives for Bayesian data analysis. This paper focuses on the subset of this interface where all measurements come from the same row. All the BQL operations used in this paper can be reduced to explicit invocations of *simulate*, *logpdf*, and to Monte Carlo estimates of Kullback-Leibler divergences implemented in terms them. Some GPMs can significantly optimize some of these operations relative to Monte Carlo baselines; such optimizations are likely to be important in practice but are beyond the scope of this paper.

### 3.1.2 WEIGHTED COLLECTIONS OF GENERATIVE POPULATION MODELS.

BQL is executed against a weighted collection of GPMs $\mathcal{M}$:

$$\mathcal{M} = \{(w_i, \mathcal{G}_i)\}$$

In principle, these collections can include GPMs drawn from different model classes. The weights are treated as prior probabilities. This paper focuses on the case where the GPMs come from a single meta-model, each produced by independent runs of a single Markov chain for posterior inference in the meta-model given all available measurements. In this case, assigning unit weights to all models $w_i = 1$ results in BQL queries based on a Monte Carlo approximation to Bayesian model averaging.

## 3.2 Core instructions: `SIMULATE`, `ESTIMATE`, and `INFER`

Data analysis workflows in BQL are built around three core classes of statistical operations:

1. Generating samples from predictive probability distributions, including both completions of existing rows in a data table as well as predictive distributions over hypothetical rows.

2. Estimating predictive probability densities and approximating derived information-theoretic quantities.

3. Summarizing multi-modal probability distributions with single values.

These capabilities are exposed via three basic extensions to SQL that each combine results from individual GPMs in different ways. They can be composed with ordinary SQL to solve a broad range of data analysis tasks:

1. *Detecting predictive relationships between variables*: `ESTIMATE COLUMN PROBABILITY OF DEPENDENCE WITH ...`

   This yields an estimate of the marginal probability of dependence between the specified columns. This is equivalent to the probability that the mutual information between those two variables is nonzero, integrating over the weighted collection of GPMs that BayesDB maintains. If the GPMs are produced by an asymptotically consistent estimator of the joint distribution, then these probabilities will reflect non-linear, heteroscedastic, or context-specific dependencies that statistical aggregates (such as correlation or linear regression coefficients) will not.

2. *Regression, classification, semi-supervised learning, and imputation*: `INFER ...`

   Each of these predictive modeling tasks requires filling in point estimates in different conditions. All of these can be viewed as special cases of `INFER`, which handles arbitrary patterns of missing values and both continuous and discrete prediction targets.

3. *Anomaly/outlier detection*: `ORDER BY PROBABILITY OF col ASCENDING LIMIT k`

   Anomalous cells can be found by predictive checking: identify the cells that are least likely under the inferred constellation of models. These may not be outliers in the standard univariate sense: the low probability may be due to interactions between several variables, even though each variable on its own is marginally typical.

4. *Retrieving similar rows*: `ORDER BY SIMILARITY TO row`

   A broad class of structured search operations can be performed via information-theoretic measures of similarity between rows. These are useful in both data exploration and in more targeted search.

5. *Predictive model checking*: `SIMULATE ...`

   By comparing aggregates from the output of `SIMULATE` to the output of the analogous `SELECT` statements, it is possible to do predictive checking without having to mention models, parameters, priors, or posteriors.

### 3.2.1 SIMULATE: GENERATING SAMPLES FROM ARBITRARY PREDICTIVE DISTRIBUTIONS.

The first, called SIMULATE, provides a flexible interface to sampling from posterior predictive distributions:

> SIMULATE *target columns* FROM *population* [WHERE *row filter*] [ASSUMING *constraint*] [*k* TIMES]

The WHERE clause is interpreted as a constraint to test against all members of the population that have been observed so far. If it is not supplied, the SIMULATE command is executed against an arbitrary as-yet-unobserved member of the population, i.e. a unique row id from the standpoint of the GPM interface. The ASSUME clause is interpreted as an additional set of constraints to condition each row on before generating the simulations.

For example, to generate a proxy dataset of two variables varA and varB, one can write SIMULATE varA, varB FROM population 100 TIMES. As another example, consider the BQL command SIMULATE varA, varD FROM population 20 TIMES WHERE varB = True AND varC IS MISSING ASSUMING varC = 3.4. This generates 20 simulated values from $p(\texttt{varA}, \texttt{varD}|\texttt{varC} = 3.4)$ for each member of the population where varB is equal to True and varC is missing. This behavior may seem non-intuitive. For example, a SIMULATE invocation with WHERE true returns $Rk$ rows, where $R$ is the number of rows in the database and $k$ is the number of output samples specified with the query. On the other hand, WHERE false yields an empty result set, always. However, this semantics allows SQL aggregates to reduce the predictions for individual source rows by grouping on the row identifiers.

To formally describe the meaning of simulate, we first introduce some notation. Let $w(\{x_{(r,c)}|c \in \mathcal{G}\})$ be the predicate denoted by the WHERE clause, i.e. $w(\cdot) = 1$ if the predicate is satisfied and 0 otherwise. Let $R$ be the set of rows for which there is at least one measurement, i.e. $R = \{r_i|(r_i, \cdot, \cdot) \in O\}$, and let $W = \{r_i|w(\{x_{(r_i,c)}|c \in \mathcal{G}\}) = 1\}$ be the set of rows that satisfy the WHERE clause's filter. If a WHERE clause is not provided, then $w(\{x_{(r,c)}|c \in \mathcal{G}\}) = 0$ for all existing rows $r \in R$, and $W = \{r^*\}$ be a set containing a single distinguished row about which no measurements are known. Let $T = \{c_i\}$ be the set of target columns, and let $A = (c_j, x_{(r,c_j)})$ be the set of assumed equality constraints. Also let $TA = T \cup \{c|(c, \cdot) \in A\}$ be the set of all columns referenced in the SIMULATE command.

For each $r^* \in W$, the SIMULATE primitive produces a set of $k$ returned realizations $S_{r^*} = \{s_i\}$ of the following generative process:

$$\mathcal{G}^i \sim Discrete(\{\mathcal{G}_j\}; w_j\})$$
$$s_i \sim \{X_{r^*,c}|c \in T\}|\{X_{r^*,c'} = x_{c'}|(c', x_{c'}) \in A\}, \Theta^i, \{z^i_{(r^*,c_m)}|c_m \in TA\}$$

This corresponds to choosing a GPM at random according to the probabilities given by their weights and then generating $s_i$ from the conditioned distribution in that model. If the models are equally weighted, i.e. $w_i = 1$, and if all the GPMs are drawn from their posterior distribution given the observations $p(\mathcal{G}|O)$, then this procedure implements sampling from the Bayesian posterior predictive distribution over the targets given all the observed data plus the additional constraints from the ASSUME clause:

$$p(\{X_{r^*,c}|c \in T\}|\{X_{r^*,c'} = x_{c'}|(c', x_{c'}) \in A\}, O)$$
$$\propto p(\{X_{r^*,c}|c \in T\}|\{X_{r^*,c'} = x_{c'}|(c', x_{c'}) \in A\}|\mathcal{G})p(\mathcal{G}|O)$$

### 3.2.2 ESTIMATE: APPROXIMATING POSTERIOR AVERAGES.

The second core BQL primitive, called ESTIMATE, allows clients to query the posterior expectations of stochastic functions that are defined over the rows and the columns:

ESTIMATE *target properties* FROM [COLUMNS OF] *table* [WHERE *row/col filter*]

**Row-wise estimands provided by BQL.** Consider the case where the rows are being queried, i.e. COLUMNS OF does not occur in the query. Let $P = \{f_i(x_{(r,c_i)}, \mathcal{G})\}$ be the set of properties whose values are requested. These properties can depend on observed measurements as well as latent components of the GPM. Let $w(\cdot)$ implement the WHERE clause's filter, as with SIMULATE. If a WHERE clause is not provided, then $w(\{x_{(r,c)}|c \in \mathcal{G}\}) = 0$ for all existing rows $r \in R$.

Given these definitions, each row in the output of this class of ESTIMATE invocations is defined as follows:

$$\{e_i\} \text{ with } e_i = \sum_k w_k f_i(x_{(r,c_i)}, \mathcal{G}_k)$$

The total set of returned rows is defined by the where clause:

$$\{\{e_i\}_r\} \text{ for } r \in W = \{r_i|w(\{x_{(r_i,c)}|c \in \mathcal{G}\}) = 1\}$$

1. $\log p = \text{predictive-probability}(\mathcal{G}, \texttt{row} = r, \texttt{col} = c)$

   This estimand is denoted PREDICTIVE PROBABILITY OF *col*, and applied against an implicitly specified row, thus picking out a single measurement in the population. It can be implemented by delegation to the underlying GPM:

   $$\text{predictive-probability}(\mathcal{G}, r, c)) = \text{logpdf}(\mathcal{G}, \emptyset, \{(r, c, x_{(r_j,c_j)} \text{ from } O_{\mathcal{G}})\})$$

   This can be used to identify outliers — measurements that are unlikely under their marginal distribution — as well as anomalous measurements that are marginally likely but unlikely given the other measurements for the same row.

2. $\text{sim(a,b)} = \text{generative-similarity}(\mathcal{G}, \text{context} = \{c_i\}, \text{rowA} = r^a, \text{rowB} = r^b)$

   Data analysts frequently want to retrieve rows from a table that are "statistically similar" to some pre-existing or hypothetical row. This is a key problem in data exploration. It is also useful when trying to explain surprising inference results or when trying to diagnose and repair data or inference quality issues. Many machine learning techniques treat similarity as a central primitive, and use a metric formulation of similarity as the basis for inductive generalization.

   Information theory provides appealing alternatives: measure similarity in terms of the amount of information one row contains about the values in another. This can be assessed against all variables or just against a "context" that is defined by a particular subset of variables. One approach, leading to a directional measure, is to measure the divergence of the distribution over values in one row from the distribution over values in another:

   $$Pr[\mathrm{D}_{KL}^{\mathcal{G}}(\vec{x}_{\mathcal{G},r^a}|_{\{c_i\}}||\vec{x}_{\mathcal{G},r^b}|_{\{c_i\}}) = 0]$$

**Column-wise estimands provided by BQL.** Another use of `ESTIMATE` is to query properties of the columns, via `ESTIMATE ... FROM COLUMNS OF ...`. Let $r^*$ be a distinguished row about which no measurements are known, i.e. $(r^*, \cdot, \cdot) \notin O$. Let $g_i(x_{(r^*,c)}, \mathcal{G})$ be a function of a set of measurements from a fresh row and the underlying GPM. It is then straightforward to define the set $G$ of values needed to check the `WHERE` filter, the set of columns $C_s$ that satisfy the filter, and the set $E$ of returned values containing all the target expressions for each satisfying column.

$$G = \{g_c(x_{(r^*,c)}, \{x_{(r,c)}|(r, c, \cdot) \in O\}, \mathcal{G}_k)|c \in \mathcal{G}_k\}$$
$$C_s = \{g|g \in G \text{ and } w(g) = 1\}$$
$$E = \cup_t \{g_t(x_{(r^*,c)}, \{(x_{(r,c)}|(r, c, \cdot) \in O\}, \}_k)|c \in C_s\}$$

1. $p = \text{marginal-dependence-prob}(\mathcal{G}, \text{colA} = c_i, \text{colB} = c_j)$

   This estimand characterizes the amount of evidence for the existence of a predictive relationship between the pair of variables $c_i$ and $c_j$. It is defined according to the information-theoretic definition of conditional independence:

   $$Pr[X_{(r^*,c_i)} \perp\!\!\!\perp X_{(r^*,c_j)}] = \sum_{\mathcal{G}} Pr[I(X_{(r^*,c_i)}; X_{(r^*,c_j)}) = 0|\mathcal{G}]Pr[\mathcal{G}]$$

   If each weighted GPM $\mathcal{G}_k$ is sampled approximately from some Bayesian posterior $Pr[\mathcal{G}|O]$ (and $w_k = 1$ identically), then simple Monte Carlo estimation of the marginal dependence probability yields an estimate of the *posterior* marginal dependence probability:
   $$Pr[X_{(r^*,c_i)} \perp\!\!\!\perp X_{(r^*,c_j)}|O]$$

2. $b = \text{mutual-information}(\mathcal{G}, \text{colA} = c_i, \text{colB} = c_j)$

   The mutual information between two columns can be estimated by the standard reduction to KL divergence (Cover and Thomas, 2012). This complements the marginal dependence probability, providing one measure of the strength of a dependence.

   For convenience, some of the quantities that are ordinarily accessed via `ESTIMATE` are also made available via `SELECT`.

### 3.2.3 `INFER`: SUMMARIZING DISTRIBUTIONS WITH POINT ESTIMATES.

Predictive modeling applications sometimes require access to point predictions rather than samples from predictive distributions. BQL provides these capabilities using the `INFER` primitive. The difference between `INFER` and `SELECT` is that `INFER` incorporates automatic implicit imputation from the underlying collection of GPMs, plus filtering based on user-specified confidence thresholds. For simplicity, this paper describes a simplified version with a single threshold:

```
INFER target columns FROM table [WHERE row filter] WITH CONFIDENCE confidence
level
```

This operation returns a set of measurements $\{x_{(r,c)}^{inf}\}$ where unobserved measurements are filled in with point estimates $\hat{x}_{(r,c)}$ if a prescribed confidence threshold $p(\text{conf}(X_{(r,c)} = \hat{x}_{(r,c)}) \geq q)$ is reached. More formally:

$$x_{(r,c)}^{inf} = \begin{cases} x_{(r,c)} & \text{foreach } (r, c, \cdot) \in O \\ \hat{x}_{(r,c)} & \text{foreach } (r, c, \cdot) \notin O \text{ and } p(\text{conf}(X_{(r,c)} = \hat{x}_{(r,c)}) \geq q) \\ \text{null} & \text{otherwise} \end{cases}$$

For discrete measurements, BQL implements $\text{conf}(\cdot)$ in terms of predictive probability:

$$\text{conf}(X_{(r,c)} = \hat{x}_{(r,c)}) = p(X_{(r,c)} = \hat{x}_{(r,c)}) = \sum_{\mathcal{G}} p(X_{(r,c)} = \hat{x}_{(r,c)}|\mathcal{G})p(\mathcal{G}) = \sum_{\mathcal{G}} p(X_{(r,c)} = \hat{x}_{(r,c)}|\mathcal{G})w_i$$

Optimal candidate estimates can be found by optimization, implemented via enumeration:

$$\hat{x}_{(r,c)} = \arg\max_{x} p(X_{(r,c)} = x)$$

For continuous measurements, there is no canonical definition of confidence that applies to all GPMs. Here we define $conf(X_{(r,c)} = x) = q$ as the probability that there is a useful unimodal summary of the distribution of $X_{(r,c)}$ that captures at least $100q$ percent of the predictive probability mass. This can be formalized in terms of mixture modeling. Let $\phi_l$ be the parameters of mixture component $l$; for continuous data, we will use Gaussian component models, so $\phi_l = (\mu_l, \sigma_l)$. Let $\pi_l$ be the mass associated with component $l$. We will choose $conf(\cdot)$ and $\hat{x}$ as follows:

$$\{(\phi_l, \pi_l)\} \sim p(\{(\phi_l, \pi_l)\}|\{X_{(r,c)}^k\}) \text{ for } 0 \leq k \leq K^+$$
$$l^* = \arg\max_{l} \pi_l$$
$$\hat{X}_{(r,c)} = \mu_{l^*}$$
$$conf(X_{(r,c)} = \hat{x}_{(r,c)}) = \pi_{l^*}$$

Note that this approach can recover the behavior of the chosen strategy for discrete data by using component models that place all their probability mass on single values. The current prototype implementation of BayesDB uses a standard Gibbs sampler for a Dirichlet process mixture model (Mansinghka et al., 2015; Neal, 1998; Rasmussen, 2000) to sample $\{(\phi_l, \pi_l)\}|\{X_{(r,c)}^k\}$, with $K^+ = 1000$ by default. Adjusting $K^+$ and the amount of inference done in this mixture model can yield a broad class of tradeoffs between time, accuracy, and variance; the current values are chosen for simplicity.

## 3.3 Model and data independence

Relational databases revolutionized the processing and analysis of business data by enabling a single centrally managed data base to shared by multiple applications and also shared between operational and analytic workloads. This in turn accelerated the development of high

performance and efficient databases, because the common abstraction became a target for researchers and industrial practitioners looking to build high performance system software with a broad impact. The relational model enabled sharing and infrastructure reuse because interactions with the data, queries, are expressed in a notation (most popularly SQL) that is independent of the physical representation of the data (Codd, 1970). Without this independence, physical data layout must be carefully tailored to particular workloads, specialized code written to manipulate the layout, and these data formats and access methods cannot easily be shared.

BayesDB aims to provide additional abstraction barriers that insulate clients from the statistical underpinnings of data analysis. Clients need to be able to specify data analysis steps and workflows in a notation that is independent of the models and runtime inference strategies used to implement individual primitives, and (where possible) the modeling strategies used to produce models from the original data.

Recall that the complete persistent state of a single population in BayesDB is characterized by two mathematical objects:

1. The complete set of observed measurements $O = \{(r_i, c_i, x_{(r_i,c_i)})\}$.

2. The weighted collection of GPMs $\{(w_k, \mathcal{G}_k)\}$. Note that this notation makes no commitment as to the content of the GPMs, the weights, or the procedures by which they were obtained.

The independencies provided by BayesDB can be described in terms of these objects:

1. *Physical data independence.* The notation for $O$ makes no commitment as to the physical representation of the measurements. The definitions of BQL primitives given above therefore do not depend on details of the data representation to define their values. However, as with SQL, small changes in representation may yield large changes in runtime performance.

2. *Physical model independence.* The notation for each $(w_k, \mathcal{G}_k)$ makes no commitment as to the specific probability distribution induced over the set of random measurements $X = \{X_{(r_i,c_i)}\}$. The definitions of BQL primitives given above therefore do not depend on the details of the probabilistic models used to define the random result set for each query. In principle, the mathematical properties of the models as well as their software implementation (or even implementing platform) can be changed without invalidating end user queries. However, small changes in the generative population model may yield large changes in the results of *simulate* and *logpdf*.

Databases provide other finer-grained independence properties that may have useful analogs in BayesDB. For example, let us partition the random variables induced by a given GPM into two subsets $X^A = \{X_{(r_i^a, c_i^a)}\}$ and $X^B = \{X_{(r_j^b, c_j^b)}\}$. An example of a desirable data-dependent independence property is that if $X^A|O \perp\!\!\!\perp X^B|O$ in the "true" GPM, then $Q|X^A, X^B = Q|X^A$ in any inferred models. Informally, this rests on the model-building strategy: if the model-builder recovers the correct independencies, then the independence of query results follows. This can be thought of as an analogue of logical data independence,

which stipulates that e.g. adding new features should not affect the behavior of existing applications whose results do not depend on the value of these new features. Formalizing and verifying these properties is an important challenge for future research.

## 4. Modeling with the Meta-Modeling Language

BayesDB also provides the Meta-Modeling Language (MML), a probabilistic programming language for building models of data tables. MML programs consist of *modeling tactics* that control the behavior of an automatic model-building engine. These tactics take several forms: statistical datatypes; initialization of weighted collections of random models; approximately Bayesian updating of the model collection; qualitative assertions about dependence and independence between variables; and the use of custom statistical models for specific conditional distributions. All these tactics are currently implemented in terms of a unifying semi-parametric Bayesian model that fills in all unspecified aspects.

### 4.1 Statistical datatypes

This metadata constrains the probability models that will be used for each column of data and can also be used to choose appropriate visualizations. It is straightforward to support several different kinds of data:

1. *Categorical values from a closed set.* This datatype includes a dictionary that maps the raw data values (often strings) to canonical numerical indexes for efficient storage and processing. This information can also be used to inform modeling tactics. For example, in the current version of MML, closed-set categorical variables are modeled generatively via a multinomial component model with a symmetric Dirichlet prior on the parameters (Mansinghka et al., 2015). Discriminative models for closed-set categorical columns could potentially use a multinomial logit link function, or an appropriate multi-class classification scheme.

2. *Binary data.* Data of this type is generatively modeled using an asymmetric Beta-Bernoulli model (Mansinghka et al., 2015) that can better handle sparse or marginally biased variables than a symmetric alternative. Also, a broad class of discriminative learning techniques can natively handle the binary classification problems induced by binary variables.

3. *Count data.* Non-negative counts can be naturally modeled generatively by a Poisson-Gamma model or discriminatively by a GLM with the appropriate link function.

4. *Numerical data.* By default, data of this type is generatively modeled using a standard Normal-Gamma model. It is straightforward to add numerical ranges to enforce truncation post-hoc, and to add numerical pre-transformations that are appropriate for data that is naturally viewed as normal only on a log scale.

We have performed preliminary experiments on other datatypes built on standard statistical models. For example, cyclic data can be handled via a von Mises model (Gopal and Yang, 2014). Many other datatypes can be handled by the appropriate generalized linear

model (McCullagh and Nelder, 1989). Broadening the set of primitive data types and assessing coverage on a representative corpus of real-world databases will be a key research challenge going forwards.

## 4.2 Bayesian generative population meta-models

Some data generators can be learned from data. Often the learning mechanism will be based on approximate probabilistic inference in a *meta-model*: a probabilistic model defined over a space of data generators, each of which is also a probabilistic model in its own right. Thus far, all BayesDB meta-models have been *Markov chain meta-models*. These meta-models internally maintain a single sample from an approximate posterior, and provide a Markov chain transition operator that updates this sample stochastically.

1. $\mathcal{G} = (\theta_{\mathcal{G}}^0, \mathbf{X}_{\mathcal{G}}) = \text{initialize}(\text{meta-schema} = \Lambda)$

   Initializes a new meta-model with arbitrary parameters and an associated tabular data store.

2. $\text{incorporate}(\text{id} = r, \text{values} = \{(c_j, x_{(r,c_j)})\})$

   Creates a new member of the population with the given row index and values and stores it. Errors result from duplicate indexes or variables $c_j$ whose values $x_{(r,c_j)}$ are not compatible with the meta-schema $\Lambda$ (e.g. because the expected data type is incompatible with a provided value).

3. $\text{remove}(\text{id} = r)$

   Removes a member of the population from the data store.

4. $\text{infer}(\text{program} = \mathcal{P})$

   Simulate an internal Markov chain transition operator $\mathcal{T}$ to improve the quality of the current sampled model representation:

$$\theta_{\mathcal{G}}^{i+1} = \mathcal{T}(\theta_{\mathcal{G}}^i)$$

Some Markov chain meta-models are *asymptotically Bayesian*, i.e. the distribution that results from sequences of $T$ updates converges to the posterior over meta-models as $T$ goes to infinity:

$$\lim_{t \to \infty} \mathrm{D}_{KL}(p(\theta_{\mathcal{G}}|\mathbf{X}_{\mathcal{G}})||p(\mathcal{T}^t(\theta_{\mathcal{G}}))) = 0$$

A sufficiently expressive Markov chain meta-model may also be asymptotically consistent in the usual sense. The default semi-parametric GPM provided by BayesDB is designed to be both asymptotically consistent and asymptotically Bayesian; these invariants are crucial for its robustness, broad applicability, and suitability for use by non-experts. Formally specifying and validating these properties is an important challenge for future research.

### 4.2.1 Controlling inference via INITIALIZE and ANALYZE.

The MML allows users to control the process by which models are created and updated to reflect the data. These capabilities are exposed via two commands:

1. `INITIALIZE k MODELS FOR population`

   This command creates models by sampling their structure and parameters from the underlying GPM's prior. This is implemented by delegation to `initialize(Λ)` where $Λ$ is the entire MML schema so far.

2. `ANALYZE [variable subset OF] population FOR timelimit`

   This command performs approximately Bayesian updating of the models in the weighted collection by delegating to the `infer()` procedure from the underlying GPM. Here is a typical invocation:

   ```
   ANALYZE my_population FOR 10 MINUTES
   ```

   When no variable subset is provided, analysis is done on all the latent variables associated with every GPM in the weighted collection. Finer-grained control is also possible using variable subset specifiers that pick out particular portions of the latent state in the GPM; these details are beyond the scope of this paper.

### 4.3 Qualitative constraints

The BayesDB MML provides constructs for specifying qualitative constraints on the dependence and independence relationships (Pearl, 1988). The model-building engine attempts to enforce them in all GPMs[1]. These constraints are specified as follows:

```
ALTER METAMODEL FOR population ENSURE colA IS [NOT] MARGINALLY DEPENDENT
ON colB
```

It is also possible to `INITIALIZE` and `ANALYZE` models that do not respect the constraints, and then enforce them after the fact:

```
ALTER MODELS FOR population ENSURE colA IS [NOT] MARGINALLY DEPENDENT
ON colB
```

These commands enable domain experts to apply qualitative knowledge to make better use of sparse data. This can be crucial for improving analysis and model credibility in the eyes of domain experts. They also create the opportunity for false or unjustified knowledge to influence the results of analysis. This can reduce credibility in the eyes of statisticians or domain skeptics who want to see all assumptions in an analysis scrutinized empirically.

---

1. The current implementation does not attempt to detect contradictions.

### 4.4 Incorporating foreign statistical models

A crucial aspect of MML is that it permits experts to override the automatic model-building machinery using custom-built statistical models. Feedforward networks of such models can be specified as follows:

```
ALTER SCHEMA FOR population MODEL output variables GIVEN input variables
USING FOREIGN PREDICTOR FROM source file
```

Presently these models are presumed to be discriminative. They are only required to be able to simulate from a probability distribution over the output variables conditioned on the inputs, and to evaluate the probability density induced by this distribution.

### 4.5 A semi-parametric factorial mixture GPM

The current implementation of MML implements all the above commands in terms of approximate inference in single, unusually flexible, semi-parametric Bayesian meta-model. This GPM is closely related to CrossCat (Mansinghka et al., 2015). The CrossCat model is a factorial Dirichlet process mixture model, where variables are assigned to specific Dirichlet process mixtures by inference in another Dirichlet process mixture model over the columns. The version used for implementing MML adds two key components:

1. *Deterministic constraints on model structure.* Users can specify constraints on the marginal dependence or independence of arbitrary pairs of variables.

2. *Feedforward networks of discriminative models conditioned on the outputs of the generative model.* This allows users to combine general-purpose density estimation with standard statistical techniques such as regression as well as complex computational models with noisy outputs.

Thus in MML, the CrossCat probability model is used as the root node in a directed graphical model. Each other node in the graph corresponds to specific discriminative model, directly conditioned on the inputs of its immediate ancestors. Undirected terms attached to the root node enforce deterministic constraints.

It is helpful to view this model in terms of a "divide and conquer" modeling strategy that bottoms out in foreign predictors and other standard parametric models from Bayesian statistics:

1. All variables not explicitly assigned to a custom model are divided into marginally independent groups. Variables in the same group are assumed to be marginally dependent. Partitions of variables that do not respect the given marginal dependence and independence constraints are rejected. Each group of variables induces an independent subproblem that will typically be far lower dimensional than the original high-dimensional problem.

2. For each subproblem, divide the rows into clusters whose values are marginally dependent given any variable-specific hyperparameters.

3. For each cluster, use a simple product of parametric models — i.e. a "naive Bayes" approach (Duda et al., 2001) — to estimate the joint distribution.

Inference in the GPM thus addresses modeling tradeoffs that resemble the decisions faced in exploratory analysis, confirmatory analysis, and predictive modeling. The most crucial decisions involve defining which subset of the data is relevant for answering each question. A secondary issue is what probabilistic model to use for each subset; absent prior knowledge, these are chosen generically, based on the type of the data in the column.

### 4.5.1 A "divide-and-conquer" generative process

The generative process that induces the default GPM can be described using the following notation:

| Name | Description |
|---|---|
| $\alpha_D$ | Concentration hyperparameter for CRP that slices the columns |
| $\vec{\lambda}_d$ | Hyperparameters for column $d$ (datatype-dependent) |
| $z_d$ | Slice (column partition) assigned to column $d$ |
| $\alpha_v$ | Concentration hyperparameter for CRP that clusters rows for slice $v$ |
| $y_r^v$ | Cluster assigned to row $r$ with respect to slice $v$ |
| $\vec{\theta}_c^d$ | Model parameters for column $d$ cluster $c$ (datatype-dependent) |
| $\vec{x}_{(\cdot,d)}^c$ | Values in cluster $c$ for column $d$, i.e. $\{x_{(r,d)} \mid y\vert_r^{z_d} = c\}$ |
| $u_d$ | An indicator such that $u_d = 1$ iff $d$ is modeled by a foreign predictor |
| $par(d)$ | The set of input dimensions for the foreign predictor conditionally modeling variable $d$ |
| $\vec{\phi}_d$ | Parameters for the foreign predictor conditionally modeling variable $l$ |
| $m_d(x_{(r,d)}; \vec{\phi}_d, \vec{x}_p)$ | The stochastic model for the foreign predictor used for variable $d$ (with density $m_d^{dens}(\cdot)$) with $\vec{x}_p = \{x_{(r,p)} \mid p \in par(d)\}$ |
| $\delta_m \vec{z}$ | Characteristic function enforcing marginal (in)dependence constraint $m$ |
| $V_d(\cdot)$ | A generic hyper-prior of the appropriate type for variable or dimension $d$. |
| $M_d(\cdot)$ and $L_D(\cdot)$ $\forall\, d\ s.t.\ u_d = 1$ | A datatype-appropriate parameter prior (e.g. a Beta prior for binary data, Normal-Gamma for continuous data, or Dirichlet for discrete data), and likelihood model (e.g. Bernoulli, Normal or Multinomial). |

Using this notation, the unconstrained generative process for the default meta-model can be concisely described in statistician's notation as follows:

$$\alpha_D \sim \text{Gamma}(k = 1, \theta = 1)$$

$$\vec{\lambda}_d \sim V_d(\cdot) \qquad\qquad \text{foreach } d \in \{1, \cdots, D\}$$

$$z_d \sim \text{CRP}(\{z_i \mid i \neq d\}; \alpha_D) \qquad\qquad \text{foreach } d \in \{1, \cdots, D\}$$

$$\alpha_v \sim \text{Gamma}(k = 1, \theta = 1) \qquad\qquad \text{foreach } v \in \vec{z}$$

$$y_r^v \sim \text{CRP}(\{y_i^v \mid i \neq r\}; \alpha_v) \qquad\qquad \text{foreach } v \in \vec{z} \text{ and}$$
$$r \in \{1, \cdots, R\}$$

$$\vec{\theta}_c^d \sim M_d(\cdot; \vec{\lambda}_d)$$

$$\vec{x}_{(\cdot, d)}^c = \{x_{(r,d)} \mid y_r^{z_d} = c\} \sim \prod_r L_d(\vec{\theta}_c^d) \qquad\qquad \text{if } u_d = 0$$

$$\vec{x}_{(\cdot, d)} = \{x_{(r,d)}\} \sim m_d(\vec{\phi}_d; \{x_{(r,p)} | p \in par(d)\}) \qquad\qquad \text{if } u_d = 1$$

$$c_m \sim \delta_m(\vec{z}) \qquad\qquad \text{foreach (in)dependence constraint}$$

The true generative process also must ensure that $c_m = 1$ for all of the $M$ (in)dependence constraints. This is enforced by conditioning on the event $\{c_m = 1\}$. A generative model for this constrained process can be given trivially by embedding the unconstrained generative process in the inner loop of a rejection sampler for $\{c_m\}$ (Mansinghka, 2009; Murray et al., 2009).

### 4.5.2 THE JOINT DENSITY

Here we use $\theta_{\mathcal{G}}$ to denote all the latent information in a semi-parametric GPM $\mathcal{G}$ needed to capture its dependence on the data $O$. This includes the concentration parameter $\alpha_D$ for the CRP over columns, the variable-specific hyper-parameters $\{\vec{\lambda}_d\}$, the column partition $\vec{z}$, the column-partition-specific concentration parameters $\{\alpha_v\}$ and row partition $\{\vec{y}^v\}$, and the category-specific parameters $\{\theta_c^d\}$. Note that in this section, $M_d, V_d, L_d$, and $CRP$ each represent probability density functions rather than stochastic simulators.

Given this notation, we have:

$$P(\theta_{\mathcal{G}}, O) = P(\mathbf{X}, \{\vec{\theta}_c^d\}, \{\vec{y}^v, \alpha_v\}, \{\vec{\lambda}_d\}, \vec{z}, \alpha_D)$$

$$= e^{-\alpha_D} \Big( \prod_{d \in D} V_d(\vec{\lambda}_d) \Big) \text{CRP}(\vec{z}; \alpha_D) \Big( \prod_{v \in \vec{z}} e^{-\alpha_v} \text{CRP}(\vec{y}^v; \alpha_v) \Big)$$

$$\times \Big( \prod_{v \in \vec{z}} \prod_{c \in \vec{y}^v} \prod_{d \in \{i \text{ s.t. } z_i = v\}} M_d(\vec{\theta}_c^d; \vec{\lambda}_d) \prod_{r \in c} L_d(x_{(r,d)}; \vec{\theta}_c^d) \Big) \Big( \prod_m \delta_m \vec{z} \Big)$$

$$\times \Big( \prod_{d \text{ with } u_d = 1} \prod_r m_d^{dens}(x_{(r,d)}; \vec{\phi}_d, \{x_{(r,p)} | p \in par(d)\}) \Big)$$

### 4.5.3 INFERENCE VIA SEQUENTIAL MONTE CARLO WITH GIBBS PROPOSALS AND GIBBS REJUVENATION

Inference in this meta-model is performed via a sequential Monte Carlo scheme, in which each row is incorporated incrementally, with all latent variables proposed from their conditional

distribution. Additionally, clients can control the frequency and target latent variables for rejuvenation kernels based on Gibbs sampling, turning the overall scheme into a resample-move algorithm (Andrieu et al., 2003; Smith et al., 2013). This combination enables parallel inference and estimation of marginal probabilities while allowing the bulk of the inferential work to be done via a suitable Markov chain.

1. incorporate(id = $r$, values = $\{(c_j, x_{(r,c_j)})\}$)

   Each row is incorporated via a single Gibbs step that numerically marginalizes out all the latent variables associated with the row (Smith et al., 2013; Murphy, 2002). The associated weight is the marginal probability of the measurements to be incorporated:

   $$w_i' = w_i * p(\{(c_j, x_{(r,c_j)})\}|\mathcal{G}_\rangle)$$

   This operation is linear in the number of observed cells for the record being incorporated, the number of total slices, and the maximum number of clusters associated with any slice.

2. infer(iterations = $N$, type = rows | columns | parameters | hyperparameters | foreign | resample, slice = $j$ | NA, cluster = $k$ | NA, foreign_predictor = $l$ | NA)

   This operation applies a particular transition operator, specified by the arguments, to a selected subset of the latent variables. Each invocation affects all particles in the sequential Monte Carlo scheme. By varying the `type` parameter, a client can control whether inference is performed over the row-cluster assignment variables, the column-slice assignment variables, the cluster parameters, the column-specific hyperparameters, or all latent variables associated with a specific foreign predictor. An invocation with `type = resample` applies multinomial resampling to the weighted collection of models.

   This allows for a limited form of inference programming (Mansinghka et al., 2014), as follows. By varying the `slice`, `cluster`, or `foreign_predictor` variables, clients can instruct the GPM to only perform inference on a specific subset of the latent variables. Computational effort can thus be focused on those latent variables that are most relevant for a given analysis, rather than uniformly distributed across all latent variables in the GPM. This is most useful when the queries of interest focus on a subset of the variables, or when the clusters are well-separated.

   The prototype implementation of BayesDB uses row-cluster, column-slice, cluster-parameter, and column-hyperparameter transition operators from Mansinghka et al. (2015). The only modification is that the log joint density now includes terms for enforcing each of the (in)dependence constraints, and also terms for the likelihood induced by each foreign predictor, as described above.

This interface allows clients to specify multiple MCMC, SMC and hybrid strategies for inference. The default inference program that is invoked by the `ANALYZE` command in BQL does no resampling and selects slices and clusters to do inference on via systematic scans. It thus can be thought of as an MCMC scheme with multiple parallel chains. This approach is conservative and makes it easier to assess the stability and reproducibility of inference, although it is unlikely to be the most efficient approach in some cases.

## 5. Discussion

This paper has described BayesDB, a probabilistic programming platform that allows users to directly query the probable implications of statistical data. The query language can solve statistical inference problems such as detecting predictive relationships between variables, inferring missing values, simulating probable observations, and identifying statistically similar database entries. Statisticians and domain experts can incorporate (in)dependence constraints and custom models using a qualitative language for probabilistic models. The default meta-model frees users from needing to know how to choose modeling approaches, remove records with missing values, detect outliers, or tune model parameters. The prototype implementation is suitable for analyzing complex, heterogeneous data tables with up to tens of thousands of rows and hundreds of variables.

### 5.1 Related work in probabilistic programming

Most probabilistic programming languages are intended for model specification (Goodman et al., 2008; Stan Development Team, 2015; Milch et al., 2007; Pfeffer, 2009). This is fundamentally different from BQL and MML:

1. In BQL, probabilistic models are never explicitly specified. Instead, an implicit set of models is averaged over (or sampled from) as needed.

2. With MML, users specify constraints on an algorithm for model discovery and need not explicitly select any specific models. These constraints generally do not uniquely identify the structure of the model that will ultimately be used.

In contrast, with languages such as Stan (Stan Development Team, 2015), each program corresponds to a specific probabilistic model whose structure is fixed by the program source. Tabular (Gordon et al., 2014), a probabilistic language designed for embedding into spreadsheets that applies user-specified factor graph models defined in terms of observed and latent variables to datasets represented as sub-tables, seems closest in structure to BQL. However, like BUGS and Stan, Tabular does not aim to hide the conceptual vocabulary of probabilistic modeling from its end users, and it focuses on user-specified models. Other integrations of probabilistic modeling with databases such as (Singh and Graepel, 2013) are similarly focused on sophisticated modeling but do not provide a model-independent abstraction for queries or support for general Bayesian data analysis.

It is straightforward to extend MML to allow syntactic escapes into all these languages that allow external probabilistic programs to be used as foreign predictors.

### 5.2 Related work in probabilistic databases

BayesDB takes a complementary approach to several recent projects that integrate aspects of probabilistic inference with databases. The most closely related systems are MauveDB (Deshpande and Madden, 2006) and BBQ (Deshpande et al., 2004). They provide *model-based views* that enable users to run standard SQL queries on the outputs of statistical models. These models must be explicitly specified as part of the schema. This is useful for some machine learning applications but does not address the core problems of applied

inference, such as data exploration, data cleaning, and confirmatory analysis. Both systems also use restricted model classes that can easily introduce substantial for ad-hoc predictive queries.

Other systems such as MLBase (Kraska et al., 2013) and GraphLab (Low et al., 2012) aim to simplify at-scale development and deployment of machine learning algorithms. MLBase and GraphLab host data in a distributed database environment and provide operators for scalable ML algorithms. Systems such as SimSQL (Cai et al., 2013) and its ancestor, MCDB (Jampani et al., 2008), provide SQL operators for efficient Monte Carlo sampling. In principle, several of these systems could serve as runtime platforms for optimized implementations of BQL and the MML.

### 5.2.1 UNCERTAIN DATA VERSUS UNCERTAIN INFERENCE

The database research community has proposed several probabilistic databases that aim to simplify the management and querying of data that is "uncertain" or "imprecise" (Dalvi et al., 2009). This "data uncertainty" is different from the inferential uncertainty that motivates BayesDB. Even when the data is known with certainty, it is rarely possible to uniquely identify a single model that can be used with complete certainty. Second, each probable model is likely to have uncertain implications. Extensions of BayesDB that augment GPMs with probabilistically coherent treatments of data uncertainty are an important area for future research.

## 5.3 Limitations and future work

Additional GPMs and meta-models are needed for some applications. There are specialized SQL databases that strike different tradeoffs between query latency, workload variability, and storage efficiency. Similarly, we expect that future GPMs and meta-models will strike different tradeoffs between prediction speed, prediction accuracy, statistical model capacity, and the amount of available data. In some cases, the semi-parametric meta-model presented here may be adequate in principle but producing an appropriate implementation is a significant systems research project. For example, it may be possible to build versions suitable for ad-hoc exploration of distributed databases such as Dremel (Melnik et al., 2010) or Spark (Zaharia et al., 2010). In other cases, fundamentally different model classes may be more appropriate. For example, it seems appealing to jointly model populations of web browsing sessions and web assets with low-dimensional latent space models (Stern et al., 2009).

It will be challenging to develop query planners that can handle GPMs given by arbitrary probabilistic programs. A key issue is that the full GPM interface allows for complex conditional queries over composite GPMs that may require data-dependent inference strategies. One potential approach is to specify GPMs as probabilistic programs in a language with programmable inference; currently, the only such language is VentureScript. The inference strategy needed to answer a given query could then be assembled on-demand.

BQL and MML have yet to incorporate key ideas from several significant subfields of statistics. For example, neither language has explicit support for causal relationships and arbitrary counterfactuals (Pearl, 1988, 2009, 2001). Both BQL and MML make the standard, simplistic assumption that data is missing at random. Neither BQL nor MML has native support for longitudinal or panel data or for time-series; instead, users must apply standard

workarounds or implement custom data types. A minor limitation is that hierarchical models are currently supported by merging subpopulations, retaining an indicator variable, and treating any variables unique to a given subpopulation as missing. It should instead be possible to build GPMs that jointly model subpopulations that are separately represented (and that therefore may not share the same set of observable variables). It will also be important to develop a formal semantics and cost model for both BQL and MML.

**Qualitative probabilistic programming.** BQL and MML are qualitative languages for quantitative reasoning. They make it possible for users to perform Bayesian data analysis without needing to know how to specify quantitative probabilities or model parameters. However, the set of qualitative constructs that they support is limited, and needs to be expanded. For example, in MML, it will be important to support conditional dependence constraints. These could be specified generatively, e.g. by defining a directed acyclic graph over subsets of variables, and leaving the model builder to fill in the (conditional) joint distributions over each subset of variables. In BQL, it would be interesting to explore the addition of commands for optimization and decision-theoretic choice, with objective functions specified both explicitly and implicitly. Finally, it will be interesting to explore elicitation strategies based on "programming by example". For example, users could create datasets by iteratively specifying prototypical examples and turn them into large datasets by treating each as the seed for a separate synthetic population, produced via `SIMULATE`.

### 5.4 Conclusion

Traditional databases protect consumers of data from "having to know how the data is organized in the machine" Codd (1970) and provide automated data representations and retrieval algorithms that perform well enough for a broad class of applications. Although this abstraction barrier is only imperfectly achieved, it has proved useful enough to serve as the basis of multiple generations of software and data systems. This decoupling of task specification from implementation made it possible to improve performance and reliability — of individual database indexes, and in some cases of entire database systems — without needing to notify end users. It also created a simple conceptual vocabulary and query language for data management and data processing that spread far farther than the systems programming knowledge needed to implement it.

BayesDB aims to insulate consumers of statistical inference from the concepts of modeling and statistics and provide a simple, qualitative interface for solving problems that currently seem quantitative and complex. It also allows models, analyses, and data resources to be improved independently. It is not yet clear how deeply the analogy with traditional databases will run. However, we hope that BayesDB represents a significant step towards making statistically rigorous empirical inference more credible, transparent and ubiquitous.

### References

Christophe Andrieu, Nando De Freitas, Arnaud Doucet, and Michael I Jordan. An introduction to mcmc for machine learning. *Machine learning*, 50(1-2):5–43, 2003.

Zhuhua Cai, Zografoula Vagena, Luis Perez, Subramanian Arumugam, Peter J Haas, and Christopher Jermaine. Simulation of database-valued markov chains using simsql. In

*Proceedings of the 2013 ACM SIGMOD International Conference on Management of Data*, pages 637–648. ACM, 2013.

Edgar F Codd. A relational model of data for large shared data banks. *Communications of the ACM*, 13(6):377–387, 1970.

Thomas M Cover and Joy A Thomas. *Elements of information theory*. John Wiley & Sons, 2012.

Nilesh Dalvi, Christopher Ré, and Dan Suciu. Probabilistic databases: diamonds in the dirt. *Communications of the ACM*, 52(7):86–94, 2009.

Amol Deshpande and Samuel Madden. Mauvedb: supporting model-based user views in database systems. In *Proceedings of the 2006 ACM SIGMOD international conference on Management of data*, pages 73–84. ACM, 2006.

Amol Deshpande, Carlos Guestrin, Samuel R Madden, Joseph M Hellerstein, and Wei Hong. Model-driven data acquisition in sensor networks. In *Proceedings of the Thirtieth international conference on Very large data bases-Volume 30*, pages 588–599. VLDB Endowment, 2004.

R.O. Duda, P.E. Hart, D.G. Stork, et al. *Pattern classification*, volume 2. wiley New York, 2001.

Andrew Gelman, John B. Carlin, Hal S. Stern, and Donald B. Rubin. *Bayesian Data Analysis*. Chapman and Hall, London, 1995.

Noah D. Goodman, Vikash K. Mansinghka, Daniel Roy, Keith Bonawitz, and Joshua B. Tenenbaum. Church: a language for generative models. In *Uncertainty in Artificial Intelligence*, 2008.

Siddharth Gopal and Yiming Yang. von mises-fisher clustering models. In *Proceedings of the 31st International Conference on Machine Learning (ICML-14)*, pages 154–162, 2014.

Andrew D Gordon, Thore Graepel, Nicolas Rolland, Claudio Russo, Johannes Borgstrom, and John Guiver. Tabular: a schema-driven probabilistic programming language. In *ACM SIGPLAN Notices*, volume 49, pages 321–334. ACM, 2014.

Ravi Jampani, Fei Xu, Mingxi Wu, Luis Leopoldo Perez, Christopher Jermaine, and Peter J Haas. Mcdb: a monte carlo approach to managing uncertain data. In *Proceedings of the 2008 ACM SIGMOD international conference on Management of data*, pages 687–700. ACM, 2008.

Tim Kraska, Ameet Talwalkar, John C Duchi, Rean Griffith, Michael J Franklin, and Michael I Jordan. Mlbase: A distributed machine-learning system. In *CIDR*, 2013.

Yucheng Low, Danny Bickson, Joseph Gonzalez, Carlos Guestrin, Aapo Kyrola, and Joseph M Hellerstein. Distributed graphlab: a framework for machine learning and data mining in the cloud. *Proceedings of the VLDB Endowment*, 5(8):716–727, 2012.

Vikash Mansinghka, Daniel Selsam, and Yura Perov. Venture: a higher-order probabilistic programming platform with programmable inference. *arXiv preprint arXiv:1404.0099*, 2014.

Vikash Mansinghka, Patrick Shafto, Eric Jonas, Cap Petschulat, Max Gasner, and Joshua Tenenbaum. Crosscat: A fully bayesian nonparametric method for analyzing heterogeneous, high dimensional data. *Journal of Machine Learning Research*, 2015.

Vikash Kumar Mansinghka. *Natively probabilistic computation*. PhD thesis, Massachusetts Institute of Technology, 2009.

Peter McCullagh and John A Nelder. *Generalized linear models*, volume 37. CRC press, 1989.

Sergey Melnik, Andrey Gubarev, Jing Jing Long, Geoffrey Romer, Shiva Shivakumar, Matt Tolton, and Theo Vassilakis. Dremel: interactive analysis of web-scale datasets. *Proceedings of the VLDB Endowment*, 3(1-2):330–339, 2010.

Brian Milch, Bhaskara Marthi, Stuart Russell, David Sontag, Daniel L Ong, and Andrey Kolobov. 1 blog: Probabilistic models with unknown objects. *Statistical relational learning*, page 373, 2007.

Kevin Patrick Murphy. *Dynamic bayesian networks: representation, inference and learning*. PhD thesis, University of California, Berkeley, 2002.

Iain Murray, David MacKay, and Ryan P Adams. The gaussian process density sampler. In *Advances in Neural Information Processing Systems*, pages 9–16, 2009.

R. Neal. Markov chain sampling methods for dirichlet process mixture models, 1998.

Judea Pearl. *Probabilistic reasoning in intelligent systems: Networks of plausible inference*. Morgan Kaufmann, San Mateo, California, 1988.

Judea Pearl. Bayesianism and causality, or, why i am only a half-bayesian. In *Foundations of bayesianism*, pages 19–36. Springer, 2001.

Judea Pearl. *Causality*. Cambridge university press, 2009.

Avi Pfeffer. Figaro: An object-oriented probabilistic programming language. *Charles River Analytics Technical Report*, 137, 2009.

C. Rasmussen. The infinite gaussian mixture model. In *Advances in Neural Processing Systems 12*, 2000.

Stuart Russell and Peter Norvig. *Artificial intelligence: A modern approach*. Prentice Hall, Upper Saddle River, New Jersey, 2003.

Sameer Singh and Thore Graepel. Automated probabilistic modeling for relational data. In *Proceedings of the ACM of Information and Knowledge Management CIKM 2013*. ACM, 2013. URL http://research.microsoft.com/apps/pubs/default.aspx?id=200220.

Adrian Smith, Arnaud Doucet, Nando de Freitas, and Neil Gordon. *Sequential Monte Carlo methods in practice.* Springer Science & Business Media, 2013.

Stan Development Team. *Stan Modeling Language Users Guide and Reference Manual, Version 2.8.0*, 2015. URL `http://mc-stan.org/`.

David H Stern, Ralf Herbrich, and Thore Graepel. Matchbox: large scale online bayesian recommendations. In *Proceedings of the 18th international conference on World wide web*, pages 111–120. ACM, 2009.

L. Wasserman. Low assumptions, high dimensions. *Rationality, Markets and Morals*, 2, 2011.

Matei Zaharia, Mosharaf Chowdhury, Michael J Franklin, Scott Shenker, and Ion Stoica. Spark: cluster computing with working sets. In *Proceedings of the 2nd USENIX conference on Hot topics in cloud computing*, volume 10, page 10, 2010.