

Personalized Health and Africa

Neil D. Lawrence

University of Sheffield

18th June 2015

Outline

Diversity of Data

Massively Missing Data

Not the Scale it's the Diversity

The screenshot shows a web browser window with the URL `dataconomy.com/big-data-proving-to-be-a-real-challenge-for-data-scientists/`. The page features the Dataconomy logo and navigation menu (NEWS, EVENTS, OPINION, START UPS, INDUSTRY, RESOURCES, ABOUT, JOBS). The article title is "Big Data Proving to Be A Real Challenge for Data Scientists" by Furhaad Shah, dated July 2, 2014. The main image is a silhouette of a person looking at a starry sky with a circular data visualization overlay. The article text discusses the challenge of diverse data types rather than just volume, quoting Marilyn Matz, CEO of Paradigm4. A quote at the bottom reads: "The increasing variety of data sources is forcing data scientists into shortcuts that leave data and money on the table," said Marilyn Matz, CEO of Paradigm4. "The focus on the volume of data hides the real challenge of data analytics today. Only diverse types of data will we be able to unlock the enormous potential of analytics."

Category: Data Science, News, [permalink](#)

Tagged under: Big Data, Data Scientist, survey

[in](#) [twitter](#) [facebook](#)

Top Stories

- Predicting the World Cup with Big Data
- Kreditech Raises \$40 Million at \$190 Million Valuation

[Privacy & Cookies Policy](#)

Outline

Diversity of Data

Massively Missing Data

Massive Missing Data

- ▶ If missing at random it can be marginalized.
- ▶ As data sets become very large (39 million in EMIS) data becomes extremely sparse.
- ▶ Imputation becomes impractical.

Imputation

- ▶ Expectation Maximization (EM) is gold standard imputation algorithm.
- ▶ Exact EM optimizes the log likelihood.
- ▶ Approximate EM optimizes a lower bound on log likelihood.
 - ▶ e.g. variational approximations (VIBES, Infer.net).
- ▶ Convergence is *guaranteed* to a local maxima in log likelihood.

Expectation Maximization

Require: An initial guess for missing data

Expectation Maximization

Require: An initial guess for missing data
repeat

Expectation Maximization

Require: An initial guess for missing data

repeat

 Update model parameters

(M-step)

Expectation Maximization

Require: An initial guess for missing data

repeat

 Update model parameters

(M-step)

 Update guess of missing data

(E-step)

Expectation Maximization

Require: An initial guess for missing data

repeat

 Update model parameters

(M-step)

 Update guess of missing data

(E-step)

until convergence

Imputation is Impractical

- ▶ In very sparse data imputation is impractical.
- ▶ EMIS: 39 million patients, thousands of tests.
- ▶ For most people, most tests are missing.
- ▶ M-step becomes confused by poor imputation.

Direct Marginalization is the Answer

- ▶ Perhaps we need joint distribution of two test outcomes,

$$p(y_1, y_2)$$

- ▶ Obtained through marginalizing over all missing data,

$$p(y_1, y_2) = \int p(y_1, y_2, y_3, \dots, y_p) dy_3, \dots, dy_p$$

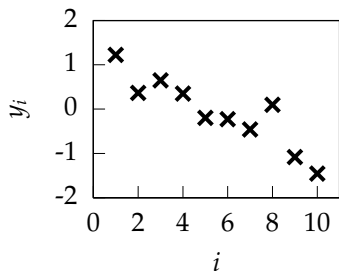
- ▶ Where y_3, \dots, y_p contains:
 1. all tests not applied to this patient
 2. all tests not yet invented!!

Magical Marginalization in Gaussians

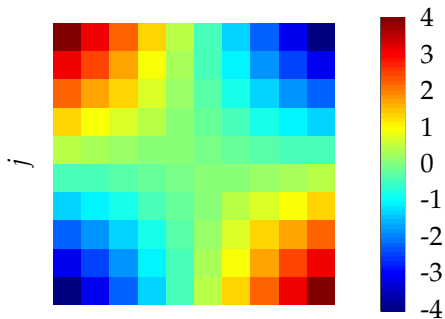
Multi-variate Gaussians

- ▶ Given 10 dimensional multivariate Gaussian, $\mathbf{y} \sim \mathcal{N}(\mathbf{0}, \mathbf{C})$.
- ▶ Generate a single correlated sample $\mathbf{y} = [y_1, y_2 \dots y_{10}]$.
- ▶ How do we find the marginal distribution of y_1, y_2 ?

Gaussian Marginalization Property



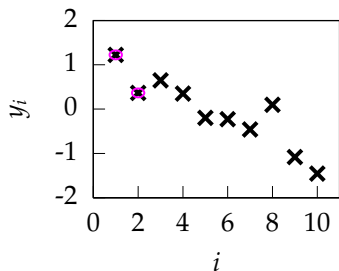
(a) A 10 dimensional sample



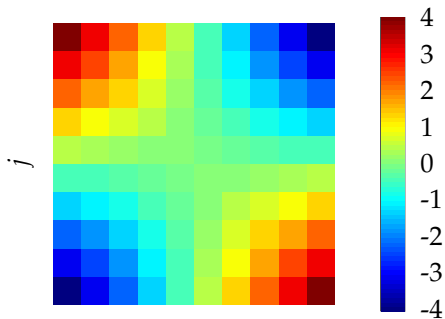
(b) colormap showing covariance between dimensions.

Figure : A sample from a 10 dimensional correlated Gaussian distribution.

Gaussian Marginalization Property



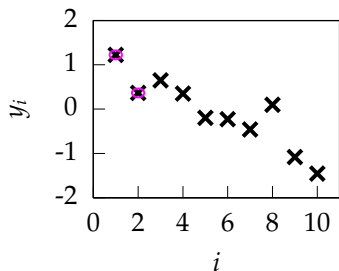
(a) A 10 dimensional sample



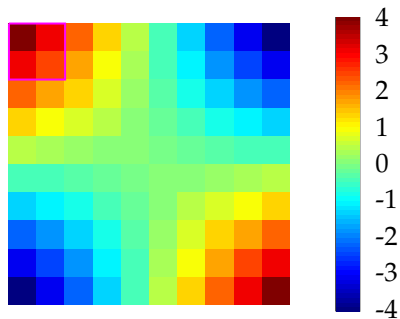
(b) colormap showing covariance between dimensions.

Figure : A sample from a 10 dimensional correlated Gaussian distribution.

Gaussian Marginalization Property



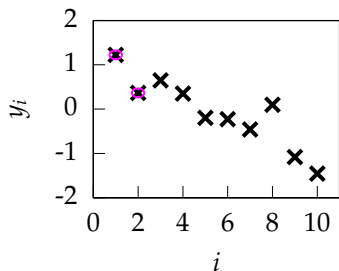
(a) A 10 dimensional sample



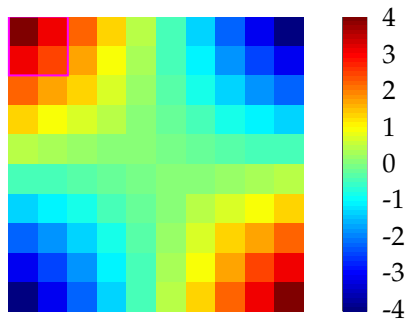
(b) colormap showing covariance between dimensions.

Figure : A sample from a 10 dimensional correlated Gaussian distribution.

Gaussian Marginalization Property



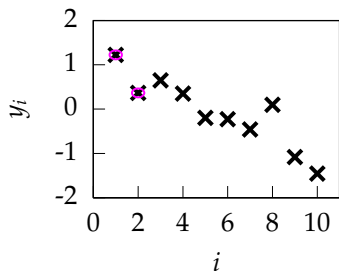
(a) A 10 dimensional sample



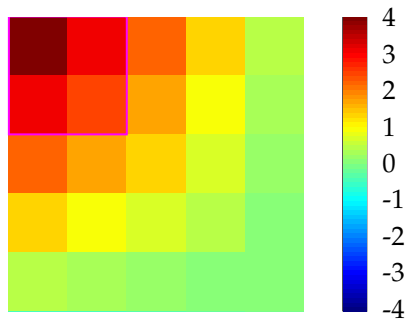
(b) colormap showing covariance between dimensions.

Figure : A sample from a 10 dimensional correlated Gaussian distribution.

Gaussian Marginalization Property



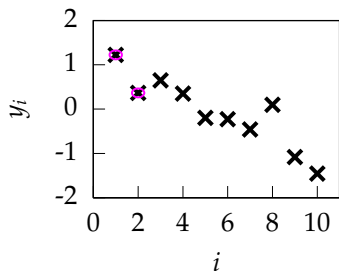
(a) A 10 dimensional sample



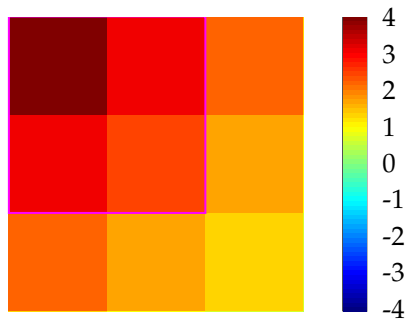
(b) colormap showing covariance between dimensions.

Figure : A sample from a 10 dimensional correlated Gaussian distribution.

Gaussian Marginalization Property



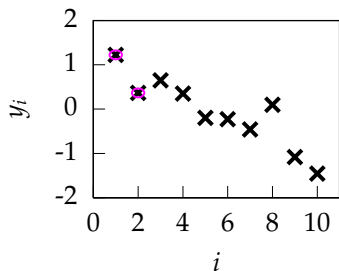
(a) A 10 dimensional sample



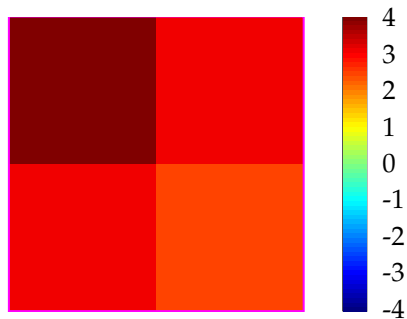
(b) colormap showing covariance between dimensions.

Figure : A sample from a 10 dimensional correlated Gaussian distribution.

Gaussian Marginalization Property



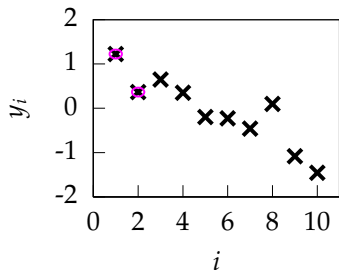
(a) A 10 dimensional sample



(b) colormap showing covariance between dimensions.

Figure : A sample from a 10 dimensional correlated Gaussian distribution.

Gaussian Marginalization Property



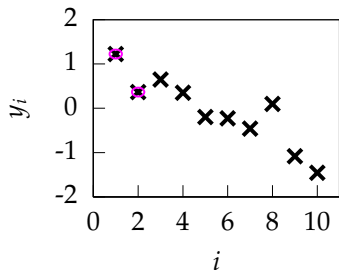
(a) A 10 dimensional sample

$$\begin{bmatrix} & 4.1 & 3.1111 \\ 3.1111 & 2.5198 & \end{bmatrix}$$

(b) covariance between y_1 and y_2 .

Figure : A sample from a 10 dimensional correlated Gaussian distribution.

Gaussian Marginalization Property



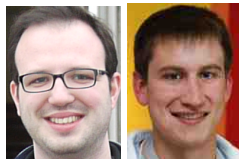
(a) A 10 dimensional sample



(b) correlation between y_1 and y_2 .

Figure : A sample from a 10 dimensional correlated Gaussian distribution.

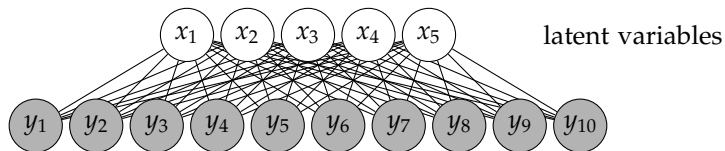
Avoid Imputation: Marginalize Directly



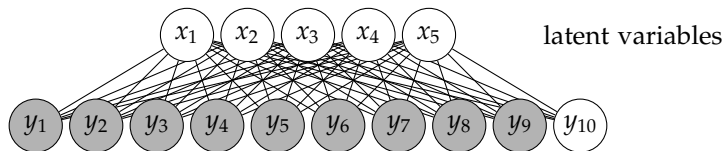
- ▶ Our approach: Avoid Imputation, Marginalize Directly.
- ▶ Explored in context of Collaborative Filtering.
- ▶ Similar challenges:
 - ▶ many users (patients),
 - ▶ many items (tests),
 - ▶ sparse data
- ▶ Implicitly marginalizes over all future tests too.

Work with Raquel Urtasun (Lawrence and Urtasun, 2009) and ongoing work with Max Zwiefsele and Nicolás Fusi.

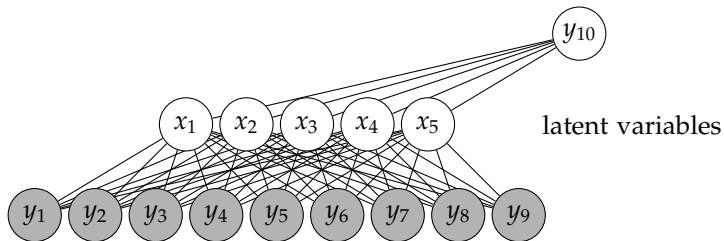
Marginalization in Bipartite Undirected Graph



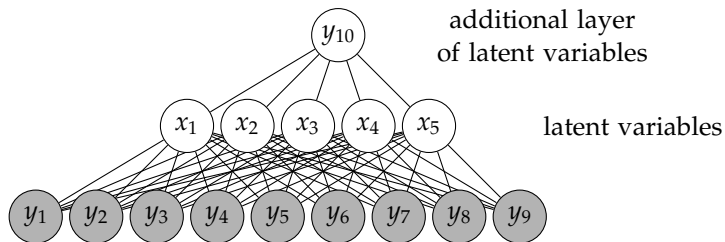
Marginalization in Bipartite Undirected Graph



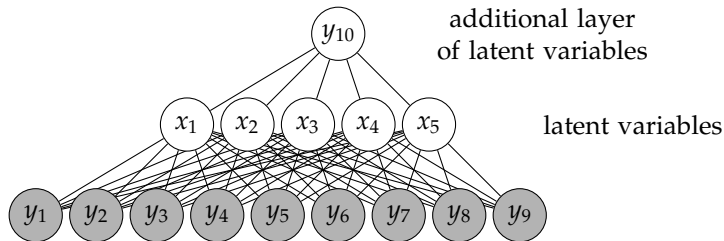
Marginalization in Bipartite Undirected Graph



Marginalization in Bipartite Undirected Graph



Marginalization in Bipartite Undirected Graph



For *massive missing data*, how many additional latent variables?

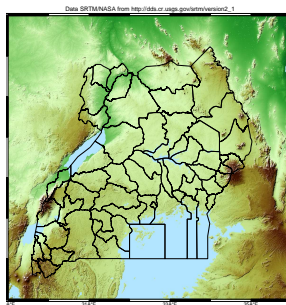
Methods that Interrelate Covariates

- ▶ Need Class of models that interrelates data, but allows for variable p .
- ▶ Common assumption: high dimensional data lies on low dimensional manifold.
- ▶ Want to retain the marginalization property of Gaussians but deal with non-Gaussian data!

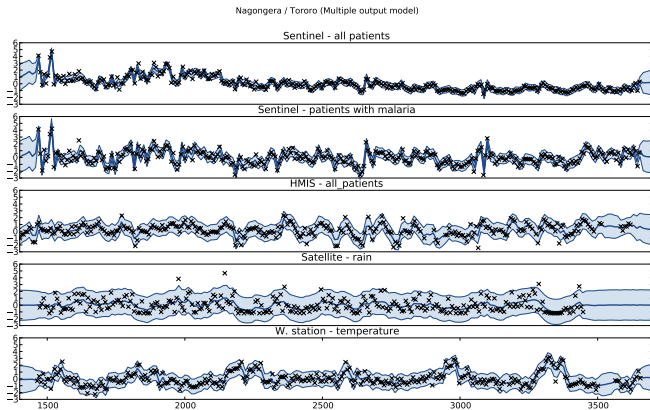
Example: Prediction of Malaria Incidence in Uganda

- ▶ Work with John Quinn and Martin Mubaganzi (Makerere University, Uganda)
- ▶ See <http://air.ug/research.html>.

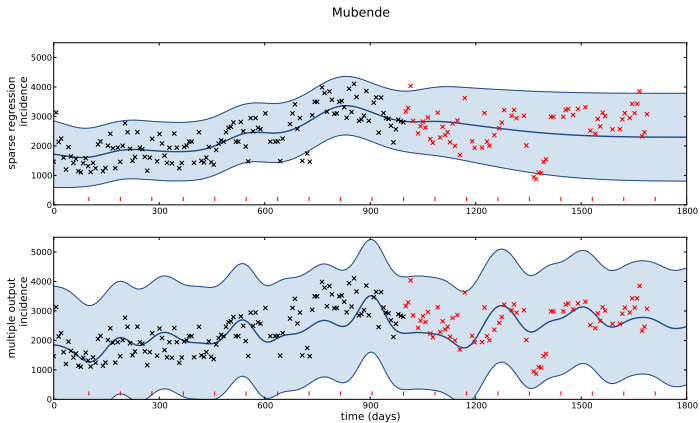
Malaria Prediction in Uganda



Malaria Prediction in Uganda



Malaria Prediction in Uganda



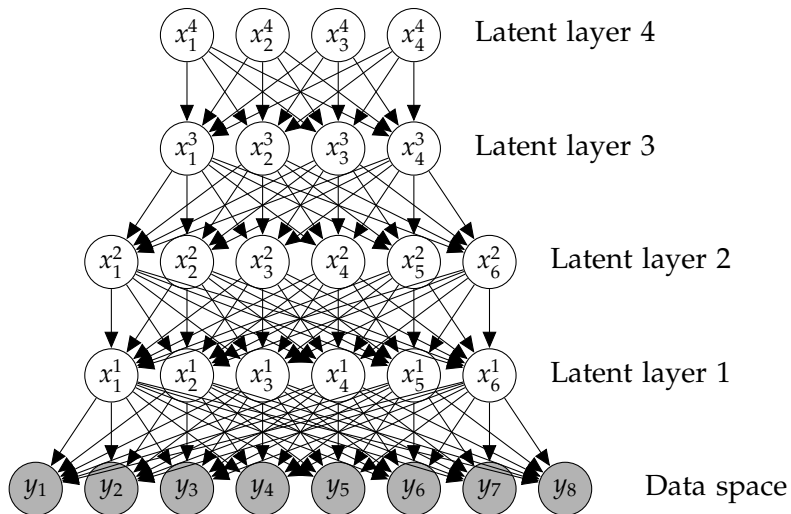
GP School at Makerere



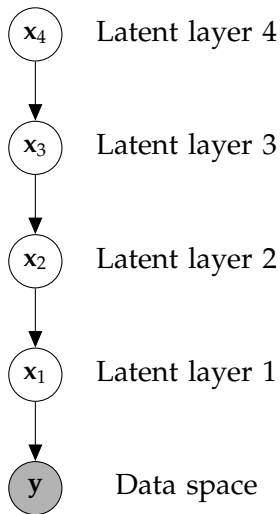
Early Warning Systems

Early Warning Systems

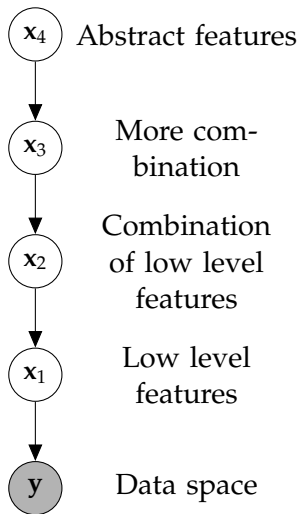
Deep Models



Deep Models



Deep Models



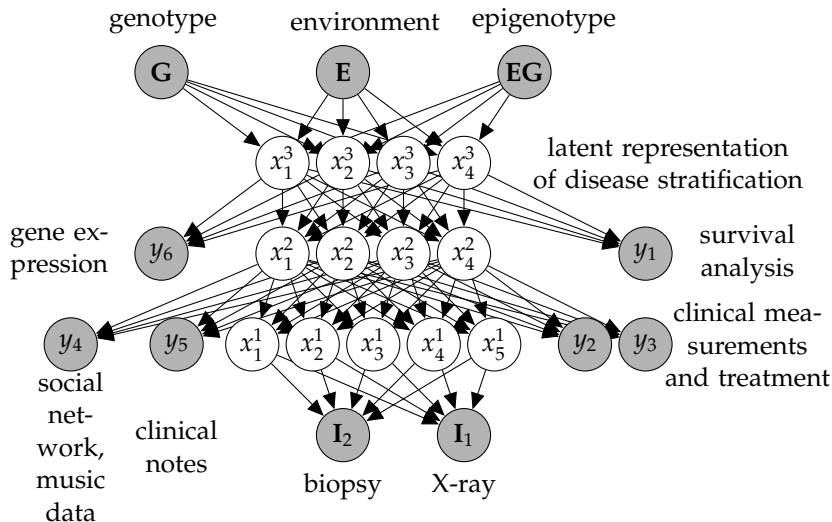
Deep Gaussian Processes



Damianou and Lawrence (2013)

- ▶ Deep architectures allow abstraction of features (Bengio, 2009; Hinton and Osindero, 2006; Salakhutdinov and Murray, 2008).
- ▶ We use variational approach to stack GP models.

Deep Health



Summary

- ▶ Intention is to deploy probabilistic machine learning for assimilating a wide range of data types in personalized health:
 - ▶ Social networking, text (clinical notes), survival times, medical imaging, phenotype, genotype, mobile phone records, music tastes, Tesco club card
- ▶ Requires population scale models with millions of features.
- ▶ May be necessary for early detection of dementia or other diseases with high noise to signal.
- ▶ Major issues in privacy and interfacing with the patient.
- ▶ But: the revolution *is* coming. We need to steer it.

References I

- Y. Bengio. Learning Deep Architectures for AI. *Found. Trends Mach. Learn.*, 2(1):1–127, Jan. 2009. ISSN 1935-8237. [\[DOI\]](#).
- A. Damianou and N. D. Lawrence. Deep Gaussian processes. In C. Carvalho and P. Ravikumar, editors, *Proceedings of the Sixteenth International Workshop on Artificial Intelligence and Statistics*, volume 31, AZ, USA, 2013. JMLR W&CP 31. [\[PDF\]](#).
- G. E. Hinton and S. Osindero. A fast learning algorithm for deep belief nets. *Neural Computation*, 18:2006, 2006.
- N. D. Lawrence and R. Urtasun. Non-linear matrix factorization with Gaussian processes. In L. Bottou and M. Littman, editors, *Proceedings of the International Conference in Machine Learning*, volume 26, San Francisco, CA, 2009. Morgan Kauffman. [\[PDF\]](#).
- R. Salakhutdinov and I. Murray. On the quantitative analysis of deep belief networks. In S. Roweis and A. McCallum, editors, *Proceedings of the International Conference in Machine Learning*, volume 25, pages 872–879. Omnipress, 2008.