



A Bayes Rule for Density Matrices

Manfred K. Warmuth

UCSC

Poster NIPS 05

Revised January 11, 2006

Thanks to Dima Kuzmin

Other Titles?

A “Quantum” Bayes Rule

or

A Non-commutative Bayes Rule

Matrix Priors

Prior probability of model M_i is $P(M_i)$

Denote prior as diagonal matrix:

$$\mathbf{I} \begin{pmatrix} \dots & & 0 \\ & P(M_i) & \\ 0 & & \dots \end{pmatrix} \mathbf{I}^T = \sum_i P(M_i) \mathbf{e}_i \mathbf{e}_i^T$$
$$\mathbf{D} \begin{pmatrix} \dots & & 0 \\ & \delta_i & \\ 0 & & \dots \end{pmatrix} \mathbf{D}^T = \sum_i \delta_i \underbrace{\mathbf{d}_i \mathbf{d}_i^T}_{\text{dyad}}$$

In above generalization of the prior,
the eigenvalues (δ_i) are probability vector
and \mathbf{D} is orthogonal system of eigenvectors

Prior now a **Density Matrix**:

- Symmetric positive definite matrix of trace one

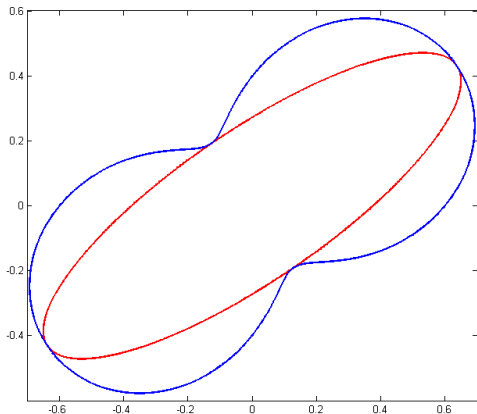
Covariance Matrices

- View the symmetric positive definite matrix \mathcal{S} as a covariance matrix of some random cost vector $\mathbf{c} \in \mathbb{R}^n$, i.e.

$$\mathcal{S} = \mathbb{E} \left((\mathbf{c} - \mathbb{E}(\mathbf{c}))(\mathbf{c} - \mathbb{E}(\mathbf{c}))^\top \right)$$

- The variance along any vector \mathbf{u} has the form

$$\mathbb{V}(\mathbf{c}^\top \mathbf{u}) = \mathbb{E} \left((\mathbf{c}^\top \mathbf{u} - \mathbb{E}(\mathbf{c}^\top \mathbf{u}))^2 \right) = \mathbb{E} \left(((\mathbf{c}^\top - \mathbb{E}(\mathbf{c}^\top)) \mathbf{u})^2 \right) = \mathbf{u}^\top \mathcal{S} \mathbf{u}$$



An ellipse \mathcal{S} in \mathbb{R}^2 : The eigenvectors are the directions of the axes and the eigenvalues their lengths. Ellipses are weighted combinations of the one-dimensional degenerate ellipses (dyads) corresponding to the axes.

The curve of the ellipse is a plot of direction \mathbf{u} times $\|\mathcal{S}\mathbf{u}\|_2$ and the outer figure eight is direction \mathbf{u} times the variance $\mathbf{u}^\top \mathcal{S} \mathbf{u}$. At the eigenvectors, this variance equals the eigenvalues and touches the ellipse.

Two More Interpretations of $\mathbf{u}^T \mathcal{S} \mathbf{u}$

- $\mathbf{u}^T \mathcal{S} \mathbf{u}$ is the squared length of \mathbf{u} w.r.t. the basis $\sqrt{\mathcal{S}}$:

$$\mathbf{u}^T \mathcal{S} \mathbf{u} = \mathbf{u}^T \sqrt{\mathcal{S}} \sqrt{\mathcal{S}} \mathbf{u} = \|\sqrt{\mathcal{S}} \mathbf{u}\|_2^2$$

- Quantum physics interpretation:

Since the squared length of \mathbf{u} w.r.t. any orthogonal basis \mathcal{S} is one, any such basis turns the unit vector into an n -dimensional probability vector $((\mathbf{u}^T \mathbf{s}_i)^2)$

Now $\mathbf{u}^T \mathcal{S} \mathbf{u}$ is the expected eigenvalue w.r.t. this probability vector:

$$\begin{aligned} \mathbf{u}^T \mathcal{S} \mathbf{u} &= \mathbf{u}^T \left(\sum_i \sigma_i \mathbf{s}_i \mathbf{s}_i^T \right) \mathbf{u} \\ &= \sum_i (\mathbf{u}^T \mathbf{s}_i)^2 \sigma_i \end{aligned}$$

Classical Setup

- Model M_i is chosen with prior probability $P(M_i)$
- Datum y is generated with probability $P(y|M_i)$

$$P(y) = \sum_i \underbrace{P(M_i)P(y|M_i)}_{\text{expected likelihood}}$$

Quantum Setup

- If prior is $\mathcal{D} = \sum_{i=1}^n \delta_i \mathbf{d}_i \mathbf{d}_i^T$, then dyad $\mathbf{d}_i \mathbf{d}_i^T$ is chosen with probability δ_i
- A random variable $\mathbf{c}^T \mathbf{d}_i$ is observed, where \mathbf{c} has covariance matrix $\mathcal{D}(y|\cdot)$
- Variance along direction \mathbf{d}_i

$$\mathbb{V}(\mathbf{c}^T \mathbf{d}_i) = \mathbf{d}_i^T \mathcal{D}(y|\cdot) \mathbf{d}_i$$

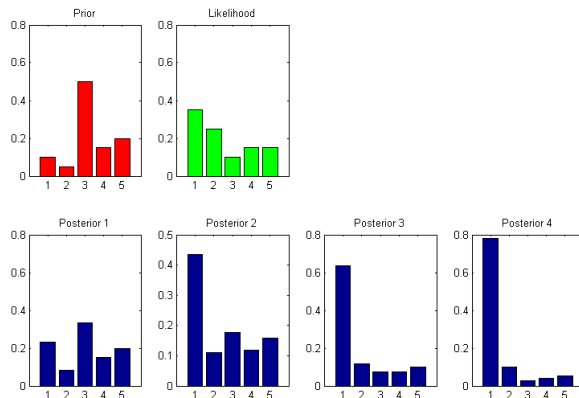


$$\underbrace{\text{tr}(\mathcal{D}(\cdot) \mathcal{D}(y|\cdot))}_{\text{quantum measurement}} = \text{tr}\left(\sum_i \delta_i \mathbf{d}_i \mathbf{d}_i^T \mathcal{D}(y|\cdot)\right) = \sum_i \delta_i \underbrace{\mathbf{d}_i^T \mathcal{D}(y|\cdot) \mathbf{d}_i}_{\text{expected variance}}$$

- When $\mathcal{D}(\cdot)$ and $\mathcal{D}(y|\cdot)$ diagonal, then in classical case

Classical Bayes Rule

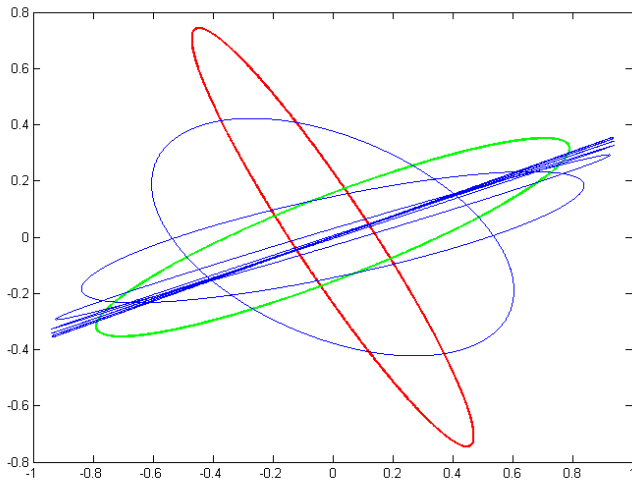
$$P(M_i|y) = \frac{P(M_i)P(y|M_i)}{P(y)}$$



- 4 updates with the same data likelihood
- Update maintains uncertainty information about maximum likelihood
- Soft max

Bayes Rule for Density Matrices

$$\mathcal{D}(\cdot|y) = \frac{\exp(\log \mathcal{D}(\cdot) + \log \mathcal{D}(y|\cdot))}{\text{tr}(\text{above matrix})}$$

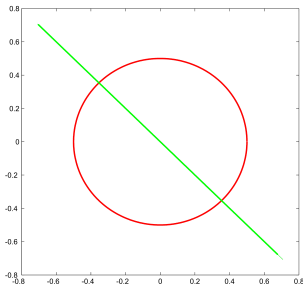


- 7 updates with same data density matrix
- Update maintains uncertainty information about maximum eigenvalue
- Soft max eigenvalue calculation

Using of Off-diagonal Elements

$$\mathcal{D}(\cdot) = \begin{pmatrix} \frac{1}{2} & 0 \\ 0 & \frac{1}{2} \end{pmatrix} \quad \mathcal{D}(y|\cdot) = \frac{1}{2} \begin{pmatrix} 1 & -1 \\ -1 & 1 \end{pmatrix} = \mathbf{U} \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix} \mathbf{U}^\top,$$

$$\text{where } \mathbf{U} = \begin{pmatrix} \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} \end{pmatrix}$$



- \forall diagonal $\mathcal{D}(\cdot)$, $\text{tr}(\mathcal{D}(\cdot) \mathcal{D}(y|\cdot)) = \frac{1}{2}$,
largest eigenvalue is not “visible” in basis \mathbf{I}

- But

$$\text{tr} \left(\underbrace{\mathbf{U} \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix} \mathbf{U}^\top}_{\mathcal{D}(\cdot|y) \text{ of new rule}} \mathcal{D}(y|\cdot) \right) = 1$$

Non-commutative Bayes Rule?

- The product of two symmetric positive matrices can be neither symmetric nor positive definite

$$\mathcal{D}(\cdot|y) = \frac{\exp\left(\log \overbrace{\mathcal{D}(\cdot)}^{\text{sym. pos. def.}} + \log \overbrace{\mathcal{D}(y|\cdot)}^{\text{sym. pos. def.}}\right)}{\text{tr}(\dots)}$$

Classical Rule as Special Case

- If $\mathcal{D}(\cdot)$ and $\mathcal{D}(y|\cdot)$ have the same eigensystem, then quantum Bayes rule specializes to the classical case

$$\begin{pmatrix} \dots & & 0 \\ & P(M_i|y) & \\ 0 & & \dots \end{pmatrix} = \frac{\begin{pmatrix} \dots & & 0 \\ & P(M_i) & \\ 0 & & \dots \end{pmatrix} \begin{pmatrix} \dots & & 0 \\ & P(y|M_i) & \\ 0 & & \dots \end{pmatrix}}{\text{tr}(\dots)}$$

Classical Intersection Properties

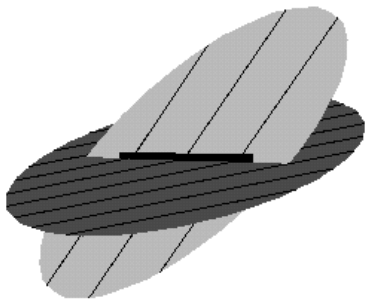
$P(M_i)$	$P(M_i y)$	$P(y M_i)$
0	0	0
a	0	0
0	b	0
a	b	ab

- Computes intersection of two sets

Generalized Case

$$\mathcal{S} \odot \mathcal{T} := \exp(\log \mathcal{S} + \log \mathcal{T}) = \lim_{n \rightarrow \infty} (\mathcal{S}^{1/n} \mathcal{T}^{1/n})^n$$

- Lie-Trotter Formula
- Limit always exists and well behaved



- “Product” lies in intersection of both spans
- In example, product is degenerate ellipse of dimension one (bold line)
- New rule:

$$\mathcal{D}(\cdot|y) = \frac{\mathcal{D}(\cdot) \odot \mathcal{D}(y|\cdot)}{\text{tr}(\mathcal{D}(\cdot) \odot \mathcal{D}(y|\cdot))}$$

Log Loss

- For classical case

$$-\log P(y) = -\log \sum_i P(M_i)P(y|M_i)$$

- In matrix case similar inequality

$$\begin{aligned} -\log \operatorname{tr}(\mathcal{D}(\cdot)\mathcal{D}(y|\cdot)) &\stackrel{\text{G.Th.}}{\leq} -\log \underbrace{\operatorname{tr}(\exp(\log \mathcal{D}(\cdot) + \log \mathcal{D}(y|\cdot)))}_{\text{normaliz. factor}} \\ &= -\log \operatorname{tr}(\mathcal{D}(\cdot) \odot \mathcal{D}(y|\cdot)) \end{aligned}$$

- Golden-Thompson inequality:

$$\forall \text{ symmetric } \mathcal{S} \text{ and } \mathcal{T} : \operatorname{tr}(\exp(\mathcal{S} + \mathcal{T})) \leq \operatorname{tr}(\exp \mathcal{S} \exp \mathcal{T})$$

- Telescoping properties of log loss still hold as inequalities
- Worst case loss bound generalize

$$\inf_w \underbrace{\Delta(\mathbf{w}, \mathbf{w}_0)}_{\text{divergence to } \mathbf{w}_0} + \underbrace{\eta}_{\text{learning rate}} \text{Loss}(\mathbf{w})$$

- Can derive large variety of updates by varying divergence, loss function and learning rate
- Examples: Gradient descent update, exponentiated gradient update, Ada-Boost ($\eta \rightarrow \infty$)
- Here we will derive Bayes rule with this framework

Classical Bayes Rule

- Mixture parameter γ_i
- Prior $P(M_i)$

$$\inf_{\gamma_i \geq 0, \sum_i \gamma_i = 1} \underbrace{\sum_i \gamma_i \log \frac{\gamma_i}{P(M_i)}}_{\text{rel.entropy}} - \eta \underbrace{\sum_i \gamma_i \log P(y|M_i)}_{\text{expected loss}}$$

- $\eta = 1$: Bayes Rule
 - Soft max
- $\eta = \infty$: maximum likelihood
- Special case of Exponentiated Gradient update

Minimization of γ

Lagrangian:

$$L(\gamma) = \sum_i \gamma_i \log \frac{\gamma_i}{P(M_i)} - \eta \sum_i \gamma_i \log P(y|M_i) + \lambda \left(\sum_i \gamma_i - 1 \right)$$

$$\frac{\partial L(\gamma)}{\partial \gamma_i} = \log \frac{\gamma_i}{P(M_i)} + 1 - \eta \log P(y|M_i) + \lambda$$

Setting partials zero:

$$\gamma_i^* = P(M_i) \exp(\lambda - 1 + \eta \log P(y|M_i))$$

Enforcing **sum constraint**:

$$\gamma_i^* = \frac{P(M_i)P(y|M_i)^\eta}{\sum_j P(M_j)P(y|M_j)^\eta}$$

$\eta = 1$: Classical Bayes rule

Bayes Rule for Density Matrices

- Parameter is density matrix \mathcal{G}
- Prior is density matrix \mathcal{D}

$$\inf_{\mathcal{G} \text{ dens.matr.}} \underbrace{\text{tr}(\mathcal{G}(\log \mathcal{G} - \log \mathcal{D}(\cdot)))}_{\text{Quantum rel. entr.}} - \eta \underbrace{\text{tr}(\mathcal{G} \log \mathcal{D}(y|\cdot))}_{\text{Fancier mixture loss}}$$

- $\eta = 1$: Quantum Bayes Rule
 - Soft maximum eigenvalue calculation
- $\eta = \infty$: minimized when \mathcal{G} is dyad uu^T and u is the eigenvector belonging to a minimum eigenvalue of

$$-\log \mathcal{D}(y|\cdot)$$

- Special case of Matrix Exponentiated Gradient update

Future

- Find conditional probability calculus for density matrices
- On-line update for PCA
- Applications
- Bayes rule for arbitrary matrices