A Bayesian Probability Calculus for Density Matrices

Manfred K. Warmuth and Dima Kuzmin

University of California - Santa Cruz

Web: Google.com "manfred"

The Technion, Haifa, Israel, Nov 4, 2013

1 Multiplicative updates - the genesis of this research

2 Conventional and Generalized Probability Distributions

3 Conventional and generalized Bayes rule

4 Bounds, derivation, calculus

1 Multiplicative updates - the genesis of this research

2 Conventional and Generalized Probability Distributions

3 Conventional and generalized Bayes rule

4 Bounds, derivation, calculus

Multiplicative updates

- A set of experts *i* learns something online
- Maintain one weight w_i per expert i
- Multiplicative update

$$w_i := \frac{w_i \ e^{-\eta \ L_i}}{\text{normalization}}$$

• Bayes rule is special case when $\eta = 1$, $L_i := -\log(P(y|M_i))$, and $w_i = P(M_i)$ $P(M_i|y) := \frac{P(M_i) P(y|M_i)}{\text{normalization}}$

- Motivated by using a relative entropy regularization
- Weight vector is uncertainty vector over experts / models

Matrix version of multiplicative updates

- Maintain density matrix W as a parameter
- Multiplicative update

$$\mathbf{W} := \frac{\exp(\log \mathbf{W} - \eta \mathbf{L})}{\text{normalization}}$$

- Motivated by using a quantum relative entropy regularization
- Density matrix is uncertainty over rank 1 subspaces
- Capping the eigenvalues at $\frac{1}{k}$: uncertainty over rank k subspaces

Matrix version of multiplicative updates

- Maintain density matrix W as a parameter
- Multiplicative update

$$\mathbf{W} := \frac{\exp(\log \mathbf{W} - \eta \mathbf{L})}{\text{normalization}}$$

- Motivated by using a quantum relative entropy regularization
- Density matrix is uncertainty over rank 1 subspaces
- Capping the eigenvalues at $\frac{1}{k}$: uncertainty over rank k subspaces
- Today: What corresponds to Bayes ???

Visualisations: ellipses

- We illustrate symmetric matrices as ellipses
 - affine transformations of the unit ball:



- Ellipse = {Su : $\|\mathbf{u}\|_2 = 1$ }
- Dotted lines connect point **u** on unit ball with point **Su** on ellipse



- For symmetric matrices, the eigenvectors form the axes of the ellipse and eigenvalues their lengths
- $\mathbf{S}\mathbf{u} = \sigma \mathbf{u}$, \mathbf{u} is an eigenvector, σ is an eigenvalue



- One eigenvalue one
- All others zero

From vectors to matrices



From vectors to matrices



i.e. mixtures of unit vectors / dyads

From vectors to matrices



Matrices are generalized distributions

Multiplicative updates - the genesis of this research

2 Conventional and Generalized Probability Distributions

3 Conventional and generalized Bayes rule

4 Bounds, derivation, calculus

Conventional probability theory

• Space is set A of n elementary events / points

 $\{a_1, a_2, a_3, a_4, a_5\}$

• Event is subset

$$\{\mathbf{a}_2, \mathbf{a}_3, \mathbf{a}_5\} = (0, \ 1, \ 1, \ 0, \ 1)^\top$$

• Distribution is probability vector

 $(.1, .2, .3, .1, .3)^{\top}$

• Probability of event $0 \cdot .1 + 1 \cdot .2 + 1 \cdot .3 + 0 \cdot .1 + 1 \cdot .3 = .8$

Conventional case in matrix notation

• The *n* elementary events are matrices with a single one on diagonal

$$egin{pmatrix} (0&0&0&0&0\\ 0&0&0&0&0\\ 0&0&0&0&0\\ 0&0&0&1&0\\ 0&0&0&0&0 \end{pmatrix} = \mathbf{e}_4 \mathbf{e}_4^ op$$

 \bullet Event $% \left\{ 0,1\right\}$ is diagonal matrix P with $\left\{ 0,1\right\}$ entries

$$\mathbf{P} = \begin{pmatrix} \begin{smallmatrix} \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \end{smallmatrix} \right) = \mathbf{e}_2 \mathbf{e}_2^\top + \mathbf{e}_3 \mathbf{e}_3^\top + \mathbf{e}_5 \mathbf{e}_5^\top$$

• Distribution is diagonal matrix **W** with probability distribution along the diagonal

$$\mathbf{W} = \begin{pmatrix} .1 & 0 & 0 & 0 & 0 \\ 0 & .2 & 0 & 0 & 0 \\ 0 & 0 & .3 & 0 & 0 \\ 0 & 0 & 0 & .1 & 0 \\ 0 & 0 & 0 & 0 & .3 \end{pmatrix}$$

Probability of event P wrt distribution W is tr(WP) = W • P

Elementary events = states in quantum physics

• Conventional: finitely many states

$$\{\mathbf{e}_i\mathbf{e}_i^ op: 1 \le i \le n\}$$

• Generalized: continuously many states $\mathbf{u}\mathbf{u}^{\top}$ called dyads

$$\{\mathbf{u}\mathbf{u}^{\top}:\mathbf{u}\in\mathbb{R}^{n},\;||\mathbf{u}||_{2}=1\}$$



- $\bullet~uu^{\bar{\top}}$ one-dimensional projection matrix onto direction u
- Prefer to use the dyads uu[⊤] as our states instead of the unit vectors u

Generalized probabilites over \mathbb{R}^n

- Elementary events are the dyads $\mathbf{u}\mathbf{u}^\top$.
- \bullet Event $% \left\{ {{\rm{B}}_{\rm{A}}} \right\}$ is symmetric matrix ${{\textbf{P}}}$ with $\left\{ {0,1} \right\}$ eigenvalues

- \mathcal{U} orthogonal eigensystem, u_i columns/eigenvectors
- **P** projection matrix onto arbitrary subspace of \mathbb{R}^n : $\mathbf{P}^2 = \mathbf{P}$

• Distribution is density matrix W:

$$\mathbf{W} = \mathcal{U} \begin{pmatrix} .1 & 0 & 0 & 0 & 0 \\ 0 & .2 & 0 & 0 & 0 \\ 0 & 0 & .3 & 0 & 0 \\ 0 & 0 & 0 & .1 & 0 \\ 0 & 0 & 0 & 0 & .3 \end{pmatrix} \mathcal{U}^{\top} = \underbrace{\sum_{i} \omega_{i} \ \mathbf{u}_{i} \mathbf{u}_{i}^{\top}}_{\text{mixture of dvads}}$$

Eigenvalues ω_i form probability vector

• Event probabilities become traces (later)

Density matrices = mixtures of dyads

• Many mixtures lead to same density matrix

$$0.2 - + 0.3 + 0.5 = \begin{pmatrix} 0.35 & 0.15 \\ 0.15 & 0.65 \end{pmatrix} = = 0.29 + 0.71$$

• There always exists a decomposition into *n* dyads that correspond to eigenvectors

Alternate defnitions of density matrices?

• Symmetric: $\mathbf{W}^{\top} = \mathbf{W}$

(

- Positive definite: $\mathbf{u}^{\top}\mathbf{W}\mathbf{u} \ge 0 \quad \forall \mathbf{u}$
- Trace one: sum of diagonal elements is one

Variance

 View the symmetric positive definite matrix C as a covariance matrix of some random cost vector c ∈ ℝⁿ, i.e.

$$\mathsf{C} = \mathbb{E}\left((\mathsf{c} - \mathbb{E}(\mathsf{c}))(\mathsf{c} - \mathbb{E}(\mathsf{c}))^{ op}
ight)$$

 $\bullet\,$ The variance along any vector u is

W

$$\begin{aligned} (\mathbf{c}^{\top}\mathbf{u}) &= \mathbb{E}(\left(\mathbf{c}^{\top}\mathbf{u} - \mathbb{E}(\mathbf{c}^{\top}\mathbf{u})\right)^{2}) \\ &= \mathbb{E}(\left((\mathbf{c}^{\top} - \mathbb{E}(\mathbf{c}^{\top}))\mathbf{u}\right)^{2}) \\ &= \mathbf{u}^{\top}\underbrace{\mathbb{E}\left((\mathbf{c} - \mathbb{E}(\mathbf{c}))(\mathbf{c} - \mathbb{E}(\mathbf{c}))^{\top}\right)}_{\mathbf{C}}\mathbf{u} \end{aligned}$$

Variance as trace

$$\mathbf{u}^{\top}\mathbf{C}\mathbf{u} = \operatorname{tr}(\mathbf{u}^{\top}\mathbf{C}\mathbf{u}) = \operatorname{tr}(\mathbf{C}\ \mathbf{u}\mathbf{u}^{\top}) = \mathbf{C} \bullet \mathbf{u}\mathbf{u}^{\top} \geq 0$$

Plotting the variance

(u^TC u)u Curve of the ellipse is plot of vector \mathbf{Cu} , where \mathbf{u} is unit vector The outer figure eight is direction **u** times the variance $\mathbf{u}^{\mathsf{T}}\mathbf{C}\mathbf{u}$ For an eigenvector, this variance equals the eigenvalue and touches the ellipse

3 dimensional variance plots



Assignment of generalized probabilities

- Density matrix W assigns generalized probability $\mathbf{u}^{\top}\mathbf{W}\mathbf{u} = \operatorname{tr}(\mathbf{W} \mathbf{u}\mathbf{u}^{\top})$ to dyad $\mathbf{u}\mathbf{u}^{\top}$
- Sum of probabilities over an orthornormal basis **u**_i is 1



For any two orthogonal directions: $\mathbf{u}_1^\top \mathbf{W} \mathbf{u}_1 + \mathbf{u}_2^\top \mathbf{W} \mathbf{u}_2 = 1$



$$a+b+c=1$$

Sum of probabilities over an orthornormal basis

۲

$$\sum_{i} \operatorname{tr}(\mathbf{W} \, \mathbf{u}_{i} \mathbf{u}_{i}^{\top}) = \operatorname{tr}(\mathbf{W} \sum_{i} \mathbf{u}_{i} \mathbf{u}_{i}^{\top}) = \operatorname{tr}(\mathbf{W} \underbrace{\mathcal{U}}_{\mathbf{I}}^{\top}) = \operatorname{tr}(\mathbf{W}) = 1$$

• Uniform density matrix: $\frac{1}{n}$



$$\operatorname{tr}(\frac{1}{n} \mathbf{I} \ \mathbf{u}\mathbf{u}^{\top}) = \frac{1}{n} \operatorname{tr}(\mathbf{u}\mathbf{u}^{\top}) = \frac{1}{n}$$

- All dyads have generalized probability $\frac{1}{n}$
- Generalized probability of *n* orthogonal dyads sums to 1

• Classical: for any probability vector $\boldsymbol{\omega}$

$$\sum_i \operatorname{tr}(\operatorname{\mathsf{diag}}(oldsymbol{\omega}) \, \mathbf{e}_i \mathbf{e}_i^ op) = \sum_i \omega_i = 1$$

Total variance along orthogonal set of directions

A density matrix For any two orthogonal directions

$$\mathbf{u}_1^{\top} \mathbf{A} \mathbf{u}_1 + \mathbf{u}_2^{\top} \mathbf{A} \mathbf{u}_2 = 1$$



$$a+b+c=1$$



Gleason's Theorem

Definition

Scalar function $\mu(\mathbf{u})$ from unit vectors \mathbf{u} in \mathbb{R}^n to \mathbb{R} is called generalized probability measure if:

• $\forall \mathbf{u}, \ \mathbf{0} \leq \mu(\mathbf{u}) \leq 1$

• If $\mathbf{u}_1, \ldots, \mathbf{u}_n$ form an orthonormal basis for \mathbb{R}^n , then $\sum \mu(\mathbf{u}_i) = 1$

Theorem

Let $n \ge 3$. Then any generalized probability measure μ on \mathbb{R}^n has the form:

$$\mu(\mathbf{u}) = \operatorname{tr}(\mathbf{W} \, \mathbf{u} \mathbf{u}^{ op})$$

for a uniquely defined density matrix W

- Disjoint events now correspond to orthogonal events
- When events have same eigensystem, then

Probability of events

• Generalized probability of event P is

$$\operatorname{tr}(\mathsf{WP}) = \operatorname{tr}(\sum_{i} \omega_{i} \ \mathbf{w}_{i} \mathbf{w}_{i}^{\mathsf{T}} \mathbf{P}) = \underbrace{\sum_{i} \omega_{i}}_{expected \ variance} \underbrace{\operatorname{variance}}_{expected \ variance}$$

 Random variable is symmetric matrix S Expectation of S is

$$\operatorname{tr}(\mathsf{WS}) = \operatorname{tr}(\mathsf{W}\sum_{i} \sigma_{i} \mathbf{s}_{i} \mathbf{s}_{i}^{\mathsf{T}}) = \underbrace{\sum_{i} \sigma_{i}}_{expected outcome} \underbrace{\mathbf{s}_{i}^{\mathsf{T}} \mathbf{W} \mathbf{s}_{i}}_{expected outcome}$$

Quantum measurement

Quantum measurement for mixture state $\mathbf{W} = \sum_{i} \omega_{i} \mathbf{w}_{i} \mathbf{w}_{i}^{\top}$ and observable $\mathbf{S} = \sum_{i} \sigma_{i} \mathbf{s}_{i} \mathbf{s}_{i}^{\top}$:

- After measurement, state collapses into one of $\{s_1s_1^{\top}, \ldots, s_ns_n^{\top}\}$
- Successor state is $\mathbf{s}_i \mathbf{s}_i^{\top}$ with probability $\mathbf{s}_i^{\top} \mathbf{W} \mathbf{s}_i$
- The expected state is again a density matrix

$$\mathbf{W} \longrightarrow \sum_{i} \mathbf{s}_{i} \mathbf{s}_{i}^{\mathsf{T}} \mathbf{W} \mathbf{s}_{i} \mathbf{s}_{i}^{\mathsf{T}} = \underbrace{\sum_{i} \mathbf{s}_{i}^{\mathsf{T}} \mathbf{W} \mathbf{s}_{i}}_{\text{expected state}} \mathbf{s}_{i} \mathbf{s}_{i} \mathbf{s}_{i}^{\mathsf{T}}$$

- Successor state $\mathbf{s}_i \mathbf{s}_i^{\top}$ associated with outcome σ_i
- Expected outcome $\operatorname{tr}(\mathbf{WS}) = \sum_i \sigma_i \mathbf{s}_i^{\top} \mathbf{W} \mathbf{s}_i$
- We use the trace computations but our density matrices are updated differently

Multiplicative updates - the genesis of this research

2 Conventional and Generalized Probability Distributions

3 Conventional and generalized Bayes rule

4 Bounds, derivation, calculus

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

$$P(y) = \sum_{i} \underbrace{P(M_{i})P(y|M_{i})}_{\text{expected likelihood}}$$



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



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



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



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



- 1 update with data likelyhood matrix
 D(y|M)
- Update maintains uncertainty information about maximum eigenvalue
- Soft max eigenvalue calculation



- 2 updates with same data likelyhood matrix D(y|M)
- Update maintains uncertainty information about maximum eigenvalue
- Soft max eigenvalue calculation


- 3 updates with same data likelyhood matrix D(y|M)
- Update maintains uncertainty information about maximum eigenvalue
- Soft max eigenvalue calculation



- 4 updates with same data likelyhood matrix D(y|M)
- Update maintains uncertainty information about maximum eigenvalue
- Soft max eigenvalue calculation



- 10 updates with same data likelyhood matrix D(y|M)
- Update maintains uncertainty information about maximum eigenvalue
- Soft max eigenvalue calculation



- 20 updates with same data likelyhood matrix D(y|M)
- Update maintains uncertainty information about maximum eigenvalue
- Soft max eigenvalue calculation

Many iterations of conventional Bayes w. same data likelihood



Plot of posterior probability as a function of the iteration number

 $(P(M_i)) = (.29, .4, .3, .01) (P(y|M_i)) = (.7, .84, .85, .9)$

Initially .85 overtakes .84 Eventually .9 overtakes both Largest likelihood: sigmoid smallest likelihood: reverse sigmoid

Many iterations of generalized Bayes Rule



Prior $\mathbf{D}(\mathbb{M})$ is diagonalized prior $(P(M_i))$ of previous plot Data likelihood $\mathbf{D}(y|\mathbb{M}) = \mathbf{U} \operatorname{diag}((P(y|M_i))\mathbf{U}^T)$,

where the eigensystem **U** is a random rotation matrix Plot of projections of the posterior onto the four eigendirections \mathbf{u}_i

Forward and backward

diagonal

with rotation





• If **D**(**M**) and **D**(**y**|**M**) have the same eigensystem, then generalized Bayes rule specializes the the conventional case

$$\begin{pmatrix} \ddots & 0 \\ P(M_i|y) \\ 0 & \ddots \end{pmatrix} = \begin{pmatrix} \ddots & 0 \\ P(M_i) \\ 0 & \ddots \end{pmatrix} \begin{pmatrix} \ddots & 0 \\ P(y|M_i) \\ 0 & \ddots \end{pmatrix}$$
$$tr(\cdots)$$

Diagonal case in dimension 2



General case in dimension 2



- Nifty generalization of Bayes rule
- Lots of evidence that it is the right generalization
- We don't have an application yet where the calculus is needed for the analysis

$$\mathsf{D}(\mathbb{M}|\mathsf{y}) = \frac{\exp\left(\mathsf{log}\,\mathsf{D}(\mathbb{M}) + \mathsf{log}\,\mathsf{D}(\mathsf{y}|\mathbb{M})\right)}{\mathrm{tr}\,(\mathsf{above\ matrix})}$$

 $\mathsf{D}(\mathbb{M})$

- symmetric positive definite
- eigenvalues sum to one

 $\textbf{D}(\textbf{y}|\mathbb{M})$

- symmetric positive definite
- all eigenvalues in range [0..1]

Later we discuss where these matrices come from

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



Conventional Bayes:

$P(M_i)$	$P(y M_i)$	$P(M_i y)$
0	0	0
а	0	0
0	b	0
а	Ь	$\frac{ab}{P(y)}$

• Computes intersection of two sets

Intersection properties on new Bayes Rule





• Result lies in intersection of both spans: here a degenerate ellipse of dimension one

Same eigensystem, then trace dot product of eigenvectors



The ω_i can track the high σ_i



- $\bullet~W$ is any matrix with eigensystem I
- \bm{S}_1, \bm{S}_2 are any matrices with rotated eigensystem and $\mathrm{tr}(\bm{S}_1) = \mathrm{tr}(\bm{S}_2),$ then

$$\operatorname{tr}(\mathbf{WS}_1) = \operatorname{tr}(\mathbf{WS}_2)$$

 \bullet So \bm{W} cannot distinguish \bm{S}_1 and \bm{S}_2

- Columns orthogonal
- $\frac{1}{\sqrt{n}}$ **H** orthonormal here called rotated eigensystem
- Dyads $\mathbf{h}\mathbf{h}^{\top}$ formed by any column \mathbf{h} of $\frac{1}{\sqrt{n}}\mathbf{H}$ has $\frac{1}{n}$ along diagonal

$$\frac{1}{n} \begin{pmatrix} 1 & -1 & 1 & -1 \\ -1 & 1 & -1 & 1 \\ 1 & -1 & 1 & -1 \\ -1 & 1 & -1 & 1 \end{pmatrix}$$

• Diagonal of
$$\mathbf{hh}^{ op}$$
 consists of all ones

Rotated versus unrotated in n dimensions



 A density matrix W with unit eigen system does not see much detail about matrices in the rotated system

Avoiding logs of zeros

Rewrite

$$\exp(\log A + \log B)$$
 as $\lim_{n \to \infty} (A^{1/n}B^{1/n})^n$

- Lie-Trotter Formula
- Limit always exists and well behaved
- Short hand for new product

$$\mathbf{A} \odot \mathbf{B} := \lim_{n \to \infty} (\mathbf{A}^{1/n} \mathbf{B}^{1/n})^n$$

• Generalized Bayes Rule becomes

$$\mathsf{D}(\mathbb{M}|\mathsf{y}) = \frac{\mathsf{D}(\mathbb{M}) \odot \mathsf{D}(\mathsf{y}|\mathbb{M})}{\operatorname{tr}(\mathsf{D}(\mathbb{M}) \odot \mathsf{D}(\mathsf{y}|\mathbb{M}))}$$

\odot - commutative matrix product

Plain matrix product is non-commutative and can violate symmetry and positive definiteness. \odot does not have these drawbacks.



Behaviour of the limit for \odot



- "Ears" indicating negative definiteness are smaller for $({\bf A}^{1/2}{\bf B}^{1/2})^2$ compared to ${\bf A}{\bf B}$
- Non-commuting part shrinks as well

Operating on bunnies with \odot



- Commutative, associative, identity matrix as neutral elmt, preserves symmetry and positive definiteness
- **2** $\mathbf{A} \odot \mathbf{B} = \mathbf{A}\mathbf{B}$ iff \mathbf{A} and \mathbf{B} commute
- **③** range(A ⊙ B) = range(A) ∩ range(B)
- $\textbf{0} \ \operatorname{tr}(\textbf{A} \odot \textbf{B}) \leq \operatorname{tr}(\textbf{A}\textbf{B}) \text{ with equality when } \textbf{A} \text{ and } \textbf{B} \text{ commute}$

● For any unit direction
$$\mathbf{u} \in \text{range}(\mathbf{A})$$
,
 $\mathbf{u}\mathbf{u}^\top \odot \mathbf{A} = e^{\mathbf{u}^\top (\log^+ \mathbf{A})\mathbf{u}} \mathbf{u}\mathbf{u}^\top$

- **(** det $(\mathbf{A} \odot \mathbf{B}) = \det(\mathbf{A}) \det(\mathbf{B})$, as for the regular matrix product
- $\textbf{O} Typically \ \textbf{A} \odot (\textbf{B} + \textbf{C}) \neq \textbf{A} \odot \textbf{B} + \textbf{A} \odot \textbf{C}$

Conventional setup again

М

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

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



Generalized setup

$\mathsf{D}(\mathsf{y}) = \operatorname{tr}(\mathsf{D}(\mathbb{M}) \odot \mathsf{D}(\mathsf{y}|\mathbb{M})) \leq \operatorname{tr}(\mathsf{D}(\mathbb{M})\mathsf{D}(\mathsf{y}|\mathbb{M}))$



where
$$\mathbf{D}(\mathbf{y}|\mathbb{M}) = \sum_i \lambda_i \mathbf{u}_i \mathbf{u}_i^{\top}$$

- Upper bounds similar to Theorem of Total Probability
- Only decouples when $D(\mathbb{M})$ and $D(y|\mathbb{M})$ have same eigensystem
 - ≤ becomes =
 - Probabilities $\rightarrow P(M_i)$ and variance/outcome $\rightarrow P(y|M_i)$

1 Multiplicative updates - the genesis of this research

2 Conventional and Generalized Probability Distributions

3 Conventional and generalized Bayes rule

4 Bounds, derivation, calculus

• Probability domain

$$P(y) = \sum_{i} P(y|M_i)P(M_i) \ge P(y|M_i)P(M_i)$$

• Log domain

$$\begin{aligned} -\log P(y) &= -\log \sum_{i} P(y|M_i) P(M_i) \\ &\leq \min_{i} (-\log P(y|M_i) - \log P(M_i)) \end{aligned}$$

Only in the log domain

• $-\log \mathbf{m}^{\top} \mathbf{A} \mathbf{m} \leq -\mathbf{m}^{\top} \log(\mathbf{A}) \mathbf{m}$, for any unit vector \mathbf{m} and symmetric positive definite matrix \mathbf{A}

$$\begin{aligned} -\log \mathbf{D}(\mathbf{y}) &= \\ &= -\log \operatorname{tr}(\mathbf{D}(\mathbf{y}|\mathbb{M}) \odot \mathbf{D}(\mathbb{M})) \\ &\leq \min_{\mathbf{m}} (-\log \mathbf{m}^{\top} (\mathbf{D}(\mathbf{y}|\mathbb{M}) \odot \mathbf{D}(\mathbb{M}))\mathbf{m}) \\ &\leq \min_{\mathbf{m}} (-\mathbf{m}^{\top} \log(\mathbf{D}(\mathbf{y}|\mathbb{M}) \odot \mathbf{D}(\mathbb{M}))\mathbf{m}) \\ &= \min_{\mathbf{m}} (-\mathbf{m}^{\top} \log \mathbf{D}(\mathbf{y}|\mathbb{M})\mathbf{m} - \mathbf{m}^{\top} \log \mathbf{D}(\mathbb{M})\mathbf{m}) \end{aligned}$$



- 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

Conventional Bayes Rule

- Mixture coefficients ω_i
- Prior $P(M_i)$

$$\inf_{\sum_{i}\omega_{i}=1} \quad \frac{1}{\eta} \underbrace{\sum_{i}\omega_{i}\log\frac{\omega_{i}}{P(M_{i})}}_{\text{rel.entropy}} - \underbrace{\sum_{i}\omega_{i}\log P(y|M_{i})}_{\text{expected loss}}$$

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

Minimization of ω

Lagrangian:

$$L(\omega) = \frac{1}{\eta} \sum_{i} \omega_{i} \log \frac{\omega_{i}}{P(M_{i})} - \sum_{i} \omega_{i} \log P(y|M_{i}) + \lambda((\sum_{i} \omega) - 1)$$
$$\frac{\partial L(\omega)}{\partial \omega_{i}} = \frac{1}{\eta} \left(\log \frac{\omega_{i}}{P(M_{i})} + 1 \right) - \log P(y|M_{i}) + \lambda$$

Setting partials zero:

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

Enforcing sum constraint:

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

 $\eta=1:$ Conventional Bayes rule



- Prior and data treated the same when $\eta = 1$
- Commutativity

Bayes Rule for density matrices

- Parameter is density matrix W
- Prior is density matrix $D(\mathbb{M})$

$$\inf_{\operatorname{tr}(\mathsf{W})=1} \quad \frac{1}{\eta} \underbrace{\operatorname{tr}(\mathsf{W}(\log \mathsf{W} - \log \mathsf{D}(\mathbb{M})))}_{\text{Quantum rel. entr.}} - \underbrace{\operatorname{tr}(\mathsf{W}\log \mathsf{D}(\mathsf{y}|\mathbb{M}))}_{\text{Fancier mixture loss}}$$

- $\eta = \infty$: minimized when **W** is dyad $\mathbf{u}\mathbf{u}^{\top}$ and **u** is the eigenvector belonging to a minimum eigenvalue of $-\log \mathbf{D}(\mathbf{y}|\mathbb{M})$
- $\eta = 1$: Generalized Bayes Rule
 - Soft maximum eigenvalue calculation
- Special case of Matrix Exponentiated Gradient update

$$\inf_{\operatorname{tr}(\mathbf{W})=1} \quad \frac{1}{\eta} \underbrace{\operatorname{tr}(\mathbf{W}(\log \mathbf{W}))}_{\text{quantum entropy}} - \frac{1}{\eta} \underbrace{\operatorname{tr}(\mathbf{W}\log \mathbf{D}(\mathbb{M}))}_{\text{initial data}} - \operatorname{tr}(\mathbf{W}\log \mathbf{D}(\mathbf{y}|\mathbb{M}))$$

- Von Neumann Entropy is just entropy of eigenvalues
- $\bullet\,$ Prior and all data treated the same when $\eta=1$
- Commutativity

- Where does data likelyhood matrix $\mathbf{D}(\mathbf{y}|\mathbb{M})$ come from?
- From a joint distribution on space (\mathbb{Y}, \mathbb{M})
Conventional joints:

- Two sets of elementary events A and B
- Joint space $A \times B$
- Elementary events are pairs (a_i, b_j)
- Joint distribution is a probability vector over pairs

Generalized joints:

- Two real vector spaces: $\mathbb A$ and $\mathbb B$ of dimension $n_{\mathbb A}$ and $n_{\mathbb B}$
- Joint space: tensor product $\mathbb{A}\otimes\mathbb{B}$ real space of dimension $n_{\mathbb{A}}n_{\mathbb{B}}$
- Elementary events are dyads of joint space
- Joint distribution is a density matrix over joint space

Joint probability?

Given joint density matrix $\mathbf{D}(\mathbb{A}, \mathbb{B})$ a dyad $\mathbf{a}\mathbf{a}^{\top}$ from space \mathbb{A} a dyad $\mathbf{b}\mathbf{b}^{\top}$ from space \mathbb{B} What's the joint probability of $\mathbf{a}\mathbf{a}^{\top}$ and $\mathbf{b}\mathbf{b}^{\top}$?

- D(a,b) =?
- Recall $\mathbf{D}(\mathbf{a}) = \operatorname{tr}(\mathbf{D}(\mathbb{A}) \mathbf{a} \mathbf{a}^{\top}).$
- Thus $\textbf{D}(a,b)=\mathrm{tr}(\textbf{D}(\mathbb{A},\mathbb{B})\ ?)$

Conventional: look up probability of jointly specified event (a_i, b_j) in joint table

What is a jointly specified dyad?

We will use Kronecker product

Kronecker product of $n \times m$ matrix **E** and $p \times q$ matrix **F** is a $np \times mq$ matrix **E** \otimes **F** which in block form is given as:

$$\mathbf{E} \otimes \mathbf{F} = \begin{pmatrix} e_{11}\mathbf{F} & e_{12}\mathbf{F} & \dots & e_{1m}\mathbf{F} \\ e_{21}\mathbf{F} & e_{22}\mathbf{F} & \dots & e_{2m}\mathbf{F} \\ \dots & \dots & \dots \\ e_{n1}\mathbf{F} & e_{n2}\mathbf{F} & \dots & e_{nm}\mathbf{F} \end{pmatrix}$$

Properties:

•
$$(\mathsf{E}\otimes\mathsf{F})(\mathsf{G}\otimes\mathsf{H})=\mathsf{E}\mathsf{G}\otimes\mathsf{F}\mathsf{H}$$

- $tr(\mathbf{E} \otimes \mathbf{F}) = tr(\mathbf{E})tr(\mathbf{F})$
- If $D(\mathbb{A})$ and $D(\mathbb{B})$ are density matrices, then so is $D(\mathbb{A}) \otimes D(\mathbb{B})$

Use $\mathbf{a}\mathbf{a}^{\top}\otimes\mathbf{b}\mathbf{b}^{\top}$ as jointly specified dyad

Joint probability: $D(a, b) = tr(D(\mathbb{A}, \mathbb{B})(aa^{\top} \otimes bb^{\top}))$

Not every dyad on the joint space can be written as $aa^{\top} \otimes bb^{\top}!!!$

This issue in quantum physics is known as entanglement

 Conditionals Marginalization Theorem of total probability

• Need additional Kronecker product properties Partial trace, etc

• Goes beyond the scope of this talk

• Many subtle quantum physics issues show up in the calculus

Sample calculus rules

•
$$\mathbf{D}(\mathbb{A}) = \operatorname{tr}_{\mathbb{B}}(\mathbf{D}(\mathbb{A}, \mathbb{B}))$$

 $\mathbf{D}(\mathbb{A}, \mathbf{b}) = \operatorname{tr}_{\mathbb{B}}(\mathbf{D}(\mathbb{A}, \mathbb{B})(\mathbf{I}_{\mathbb{A}} \otimes \mathbf{b}\mathbf{b}^{\top}))$
Marginalization

- D(A|B) = D(A, B) ⊙ (I_A ⊗ D(B))⁻¹ Conditional in terms of the joint Introduced by Cerf and Adami
- $D(\mathbb{A}) = \operatorname{tr}_{\mathbb{B}}(D(\mathbb{A}|\mathbb{B}) \odot (I_{\mathbb{A}} \otimes D(\mathbb{B})))$ Theorem of total probability
- $D(\mathbb{M}|\mathbf{y}) = \frac{D(\mathbb{M}) \odot D(\mathbf{y}|\mathbb{M})}{\operatorname{tr}(D(\mathbb{M}) \odot D(\mathbf{y}|\mathbb{M}))}$ Our Bayes rule
- $D(\mathbf{b}|\mathbb{A}) = D(\mathbf{b})D(\mathbb{A}|\mathbf{b}) \odot (D(\mathbb{A}|\mathbb{B}) \odot (I_{\mathbb{A}} \otimes D(\mathbb{B})))^{-1}$ Another Bayes rule

- We maintain uncertainty about direction of maximum variance with a density matrix
- Update generalizes conventional Bayes's rule
- Motivate the update based on a maxent principle
- Probability calculus that retains conventional probabilities as a special case

- Calculus for other matrix classes
- On-line update for PCA :-)
- Applications that need the calculus
- Connections to quantum computation

• You can implement the conventional Bayes Rule in the tube - in vitro selection with 10^{20} variable

- $\bullet\,$ You can implement the conventional Bayes Rule in the tube $\,$ in vitro selection with 10^{20} variable
- Recall Generalized Bayes Rule?

$$\mathsf{D}(\mathbb{M}|\mathsf{y}) = \frac{\mathsf{exp}\left(\mathsf{log}\,\mathsf{D}(\mathbb{M}) + \mathsf{log}\,\mathsf{D}(\mathsf{y}|\mathbb{M})\right)}{\mathrm{tr}\,(\mathsf{above\ matrix})}$$

- You can implement the conventional Bayes Rule in the tube in vitro selection with 10^{20} variable
- Recall Generalized Bayes Rule?

$$\mathsf{D}(\mathbb{M}|\mathsf{y}) = \frac{\mathsf{exp}\left(\mathsf{log}\,\mathsf{D}(\mathbb{M}) + \mathsf{log}\,\mathsf{D}(\mathsf{y}|\mathbb{M})\right)}{\mathrm{tr}\,(\mathsf{above\ matrix})}$$

• Is there a physical realization?

• The Generalized Bayes Rule retains the Conventional Bayes Rule as the hardest case

• Learning the eigensystem for free

• The Generalized Bayes Rule retains the Conventional Bayes Rule as the hardest case

• Learning the eigensystem for free

• Find other cases where there is a "free matrix lunch"