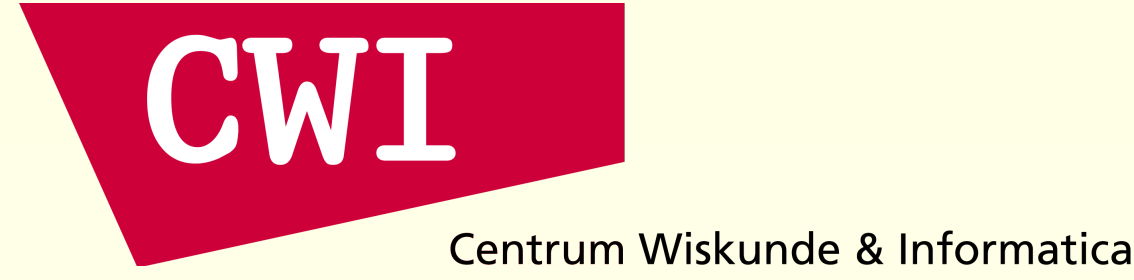


Wouter M. Koolen



Wojciech Kotłowski



Manfred K. Warmuth



## The open problem

Recent interest in matrix generalizations of classical prediction tasks:

- Matrix Hedge (for PCA)
- M. Winnow (learning subspaces)
- M. Exponentiated Gradient (regression)

In each case the matrix generalizations of classical algorithms have performance guarantees (worst-case regret bounds) *identical* to the classical tasks

Symmetric matrices have  $n^2$  parameters and vectors  $n$  parameters. Thus matrices should be *harder* to learn!

Are classical bounds loose, or is there a

Free Matrix Lunch?

## Our contribution

Extend the classical problem of predicting a sequence of outcomes from finite alphabet to matrix domain

	classical	matrix
outcomes	set of size $n$	unit vectors in $\mathbb{R}^n$
uncertainty	multinomial	density matrix
parameters	$n$	$n^2$

We show how popular online algorithms for learning a multinomial distribution can be extended to learn density matrices

Learning the  $n^2$  parameters of a density matrix should be much harder than learning the  $n$  parameters of a multinomial distribution

Surprising, we prove that the worst-case regrets of the classical algorithms and their matrix generalizations are identical:

Free Matrix Lunch!

Algorithms incur no overhead for learning the eigenvectors of the density matrix

## Many open problems

- Does the free matrix lunch hold for the matrix minimax algorithm? cf. Shtarkov
- Same questions for other losses
- What properties of the loss function and algorithm cause the free matrix lunch to occur? Proper scoring rules?
- Is there a general regret-bound preserving lift of classical algorithms to matrix prediction?

## Probability vector prediction

for trial  $t = 1, 2, \dots$  do  
 Alg predicts with probability vector  $\omega_t$   
 Nat returns outcome  $x_t$   
 Alg incurs loss  $-\log \omega_t[x_t]$   
 end for

## Density matrix prediction

for trial  $t = 1, 2, \dots$  do  
 Alg predicts with density matrix  $W_t$   
 Nat returns unit vector  $x_t$   
 Alg incurs loss  $-\mathbf{x}_t^\top \log(W_t) \mathbf{x}_t$   
 end for

## Regret

Loss of the algorithm minus the loss of the best fixed prediction in hindsight

Goal: design online algorithms with low regret

Regret of classical **Laplace** predictor  $\omega_{t+1} = \frac{\sum_{q=1}^t e_{x_q} + \mathbf{1}}{t+n}$  where  $e_i$  is  $i$ th basis vector, is

$$\text{Regret}(x_1, \dots, x_T) \leq (n-1) \log(T+1)$$

Regret of classical **Krychevsky-Trofimoff** predictor  $\omega_{t+1} = \frac{\sum_{q=1}^t e_{x_q} + \mathbf{1}/2}{t+n/2}$  is

$$\text{Regret}(x_1, \dots, x_T) \leq \frac{n-1}{2} (\log(T+1) + \log(\pi))$$

## Density matrix

Positive-semidefinite matrix  $A$  of unit trace  
 Decomposition:

$$A = \sum_i \alpha_i \mathbf{a}_i \mathbf{a}_i^\top$$

where eigenvalues  $\alpha$  form probability vector and eigenvectors  $\mathbf{a}_i$  are orthonormal system

## Log loss

The log loss is the fundamental loss for **forecasting - data compression - investment**  
 In matrix case, discrepancy between density matrix prediction  $W$  and unit vector outcome  $x$  is measured by the *matrix log loss*

$$-\mathbf{x}^\top \log(W) \mathbf{x}$$

Equal to quantum cross entropy

## Quantum entropy

$$H(A) = -\text{tr}(A \log A),$$

for density matrix  $A$  with *matrix logarithm*

$$\log A = \sum_i \log(\alpha_i) \mathbf{a}_i \mathbf{a}_i^\top$$

Equal to Shannon entropy of eigenvalues  $\alpha$

## Matrix log loss is proper

The cumulative loss of a fixed prediction  $W$

$$\sum_{t=1}^T -\mathbf{x}_t^\top \log(W) \mathbf{x}_t$$

is minimized at the empirical mean  $W^* = \frac{1}{T} \sum_{t=1}^T \mathbf{x}_t \mathbf{x}_t^\top$ , with value equal to  $T$  times the quantum entropy:  $T H(W^*)$

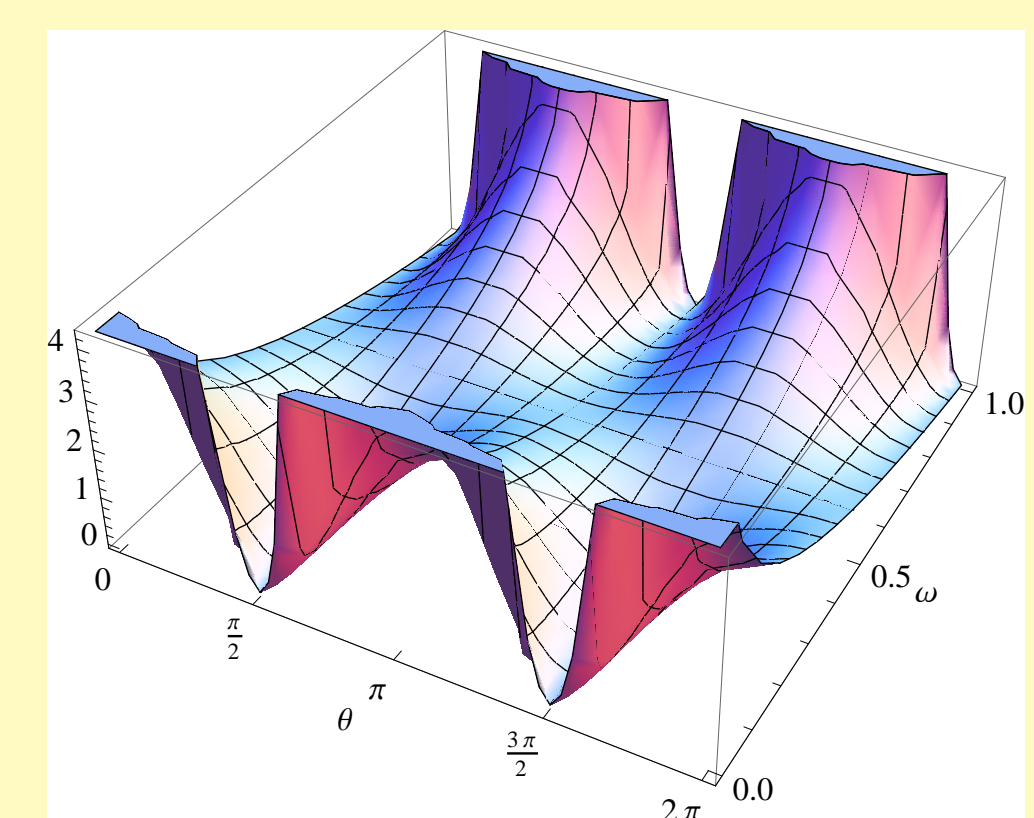
## Example: 2D matrix log loss

In  $n = 2$  dimensions, we can parametrize the prediction and outcome as follows:

$$W = \begin{pmatrix} \omega & 0 \\ 0 & 1 - \omega \end{pmatrix} \quad \text{and} \quad \mathbf{x} = \begin{pmatrix} \cos \theta \\ \sin \theta \end{pmatrix}$$

with  $\omega \in [0, 1]$  and  $\theta \in [0, 2\pi]$ . The loss becomes

$$-\mathbf{x}^\top \log(W) \mathbf{x} = -\cos^2 \theta \log \omega - \sin^2 \theta \log(1 - \omega)$$



## Result: worst-case classical and matrix regret coincide

for both **Matrix Laplace**

$$W_{t+1} = \underset{W \text{ dens. mat.}}{\text{argmin}} \left\{ \underbrace{-\text{tr}(\log W)}_{n \text{ uniform outcomes}} + \sum_{q=1}^t -\mathbf{x}_q^\top \log(W) \mathbf{x}_q \right\} = \frac{\sum_{q=1}^t \mathbf{x}_q \mathbf{x}_q^\top + \mathbf{I}}{t+n}$$

and **Matrix KT**

$$W_{t+1} = \underset{W \text{ dens. mat.}}{\text{argmin}} \left\{ \underbrace{-\frac{1}{2} \text{tr}(\log W)}_{n/2 \text{ uniform outcomes}} + \sum_{q=1}^t -\mathbf{x}_q^\top \log(W) \mathbf{x}_q \right\} = \frac{\sum_{q=1}^t \mathbf{x}_q \mathbf{x}_q^\top + \mathbf{I}/2}{t+n/2}$$

**Any sequence of outcomes not in same eigensystem is suboptimal for Nat**