Learning Eigenvectors for Free

NIPS poster Granada 12.12.2011



• 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!

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 $\omega_{t+1} = \frac{\sum_{q=1}^{t} e_{x_q} + 1}{t+n}$

where e_i is *i*th basis vector, is

Are classical bounds loose, or is there a



Our contribution

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

	classical	matrix
outcomes	set of size <i>n</i>	unit vectors in \mathbb{R}^n
uncertainty	multinomial	density matrix

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

Regret of classical **Krychevsky-Trofimoff** predictor $\boldsymbol{\omega}_{t+1} = \frac{\sum_{q=1}^{t} \boldsymbol{e}_{x_q} + 1/2}{t+n/2}$ is $\operatorname{Regret}(x_1, \dots, x_T) \leq \frac{n-1}{2} \left(\log(T+1) + \log(\pi) \right)$

Density matrix

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

$$\boldsymbol{A} = \sum_{i} \alpha_{i} \boldsymbol{a}_{i} \boldsymbol{a}_{i}^{\mathsf{T}}$$

where eigenvalues α form probability vector and eigenvectors a_i are orthonormal system

Quantum entropy

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

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*

 $-oldsymbol{x}^ op \log(oldsymbol{W})oldsymbol{x}$

Equal to quantum cross entropy

Matrix log loss is proper

The cumulative loss of a fixed prediction W

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:



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

for density matrix *A* with *matrix logarithm*

$$\log \boldsymbol{A} = \sum_{i} \log(\alpha_i) \boldsymbol{a}_i \boldsymbol{a}_i^{\top}$$

Equal to Shannon entropy of eigenvalues α

Example: 2D matrix log loss

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

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

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

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

Result: worst-case classical and matrix regret coincide

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

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



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?

for both Matrix Laplace

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

and Matrix KT
$$W_{t+1} = \operatorname*{argmin}_{W \text{ dens. mat.}} \left\{ -\frac{1}{2} \operatorname{tr}(\log W) + \sum_{q=1}^{t} -x_q^{\top} \log(W) x_q \right\} = \frac{\sum_{q=1}^{t} x_q x_q^{\top} + \mathbf{I}/2}{t+n/2}$$

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

n/2 uniform outcomes