

Learning Eigenvectors for Free

Wouter M. Koolen



Wojciech Kotłowski



Manfred K. Warmuth



WITMSE 2012

Sunday 30th September, 2012

From learning vectors to learning matrices

- Machine learning is traditionally interested in learning **vector parameters** (e.g. regression, classification)

$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} \quad \mathbf{w} = \begin{pmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{pmatrix}$$

From learning vectors to learning matrices

- Machine learning is traditionally interested in learning **vector parameters** (e.g. regression, classification)

$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} \quad \mathbf{w} = \begin{pmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{pmatrix}$$

- Recent interest in **matrix generalizations** of classical prediction tasks (PCA, learning kernels, learning subspaces)

$$\mathbf{X} = \begin{pmatrix} x_{1,1} & x_{1,2} & \cdots & x_{1,n} \\ x_{2,1} & x_{2,2} & \cdots & x_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n,1} & x_{n,2} & \cdots & x_{n,n} \end{pmatrix} \quad \mathbf{W} = \begin{pmatrix} w_{1,1} & w_{1,2} & \cdots & w_{1,n} \\ w_{2,1} & w_{2,2} & \cdots & w_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ w_{n,1} & w_{n,2} & \cdots & w_{n,n} \end{pmatrix}$$

The open problem (Warmuth, COLT 2007)

- In each case the matrix generalizations have performance guarantees (worst-case regret bounds) **identical** to the classical tasks
- Matrices have n^2 parameters and vectors n parameters. Thus matrices should be **harder** to learn!

The open problem (Warmuth, COLT 2007)

- In each case the matrix generalizations have performance guarantees (worst-case regret bounds) **identical** to the classical tasks
- Matrices have n^2 parameters and vectors n parameters. Thus matrices should be **harder** to learn!



Free Matrix Lunch???

- Predicting n -ary sequence with logarithmic loss
 - Many interpretations: forecasting, data compression, investment
 - Simple but fundamental
 - Extremely well-studied

This talk

- Predicting n -ary sequence with logarithmic loss
 - Many interpretations: forecasting, data compression, investment
 - Simple but fundamental
 - Extremely well-studied
- We generalise the problem and lift the algorithms to the matrix domain.

- Predicting n -ary sequence with logarithmic loss
 - Many interpretations: forecasting, data compression, investment
 - Simple but fundamental
 - Extremely well-studied
- We generalise the problem and lift the algorithms to the matrix domain.
- We prove and explain a



Outline

- 1 Introduction
- 2 Classical Log Loss
- 3 Matrix Log Loss
- 4 Free Matrix Lunch
- 5 Summary and Open Questions

Predicting outcomes from individual n -ary sequence (a.k.a. universal coding for n -ary alphabet)

for trial $t = 1, 2, \dots$ **do**

Alg predicts with a distribution ω_t on n -ary alphabet

Nat reveals an outcome $x_t \in \{1, \dots, n\}$

Alg incurs loss $-\log \omega_t(x_t)$

end for

Predicting outcomes from individual n -ary sequence (a.k.a. universal coding for n -ary alphabet)

for trial $t = 1, 2, \dots$ **do**

Alg predicts with probability vector (distribution) ω_t

Nat reveals a basis vector $\mathbf{x}_t \in \{e_1, \dots, e_n\}$

Alg incurs loss $-\log(\omega_t^\top \mathbf{x}_t)$

end for

Predicting outcomes from individual n -ary sequence (a.k.a. universal coding for n -ary alphabet)

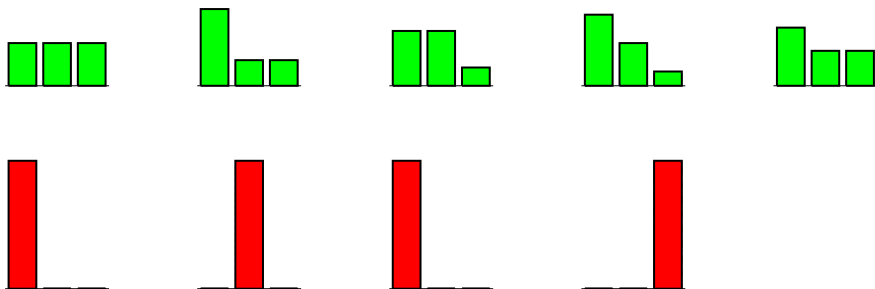
for trial $t = 1, 2, \dots$ **do**

Alg predicts with probability vector (distribution) ω_t

Nat reveals a basis vector $x_t \in \{e_1, \dots, e_n\}$

Alg incurs loss $-\log(\omega_t^\top x_t)$

end for



- **Regret** is the cumulative loss of **Alg** minus the loss of the **best fixed distribution (prediction)**:

$$\mathcal{R}_T := \sum_{t=1}^T -\log(\omega_t^\top \mathbf{x}_t) - \min_{\omega} \sum_{t=1}^T -\log(\omega^\top \mathbf{x}_t).$$

- **Regret** is the cumulative loss of **Alg** minus the loss of the **best fixed distribution (prediction)**:

$$\mathcal{R}_T := \sum_{t=1}^T -\log(\omega_t^\top \mathbf{x}_t) - \min_{\omega} \sum_{t=1}^T -\log(\omega^\top \mathbf{x}_t).$$

- The best distribution $\omega^* = \arg \min \sum_{t=1}^T -\log(\omega^\top \mathbf{x}_t)$ is the **maximum likelihood estimator**, while the loss of ω^* is the **empirical Shannon entropy**:

$$\omega^* = \frac{1}{T} \sum_{t=1}^T \mathbf{x}_t, \quad \inf_{\omega} \sum_{t=1}^T -\log(\omega^\top \mathbf{x}_t) = T H(\omega^*)$$

- **Goal**: design online algorithms with **low worst-case regret**

- Laplace predictor:

$$\omega_{t+1} := \frac{\sum_{q=1}^t \mathbf{x}_q + \mathbf{1}}{t + n} \quad \mathcal{R}_T \leq (n - 1) \log T + O(1)$$

- Laplace predictor:

$$\omega_{t+1} := \frac{\sum_{q=1}^t \mathbf{x}_q + \mathbf{1}}{t + n} \quad \mathcal{R}_T \leq (n - 1) \log T + O(1)$$

- Krychevsky-Trofimoff (KT) predictor:

$$\omega_{t+1} := \frac{\sum_{q=1}^t \mathbf{x}_q + \mathbf{1}/2}{t + n/2} \quad \mathcal{R}_T \leq \frac{n - 1}{2} \log T + O(1)$$

- Minimax regret achieved by **Shtarkov (NML)** algorithm:

$$\mathcal{R}_T \leq \frac{n-1}{2} \log T + O(1)$$

- Minimax regret achieved by **Shtarkov (NML)** algorithm:

$$\mathcal{R}_T \leq \frac{n-1}{2} \log T + O(1)$$

- **Last Step Minimax** algorithm

$$\mathcal{R}_T \leq \frac{n-1}{2} \log T + O(1)$$

Optimal up to $O(1)$. Beats KT (by a constant).

Outline

- 1 Introduction
- 2 Classical Log Loss
- 3 Matrix Log Loss**
- 4 Free Matrix Lunch
- 5 Summary and Open Questions

Density matrix prediction

for trial $t = 1, 2, \dots$ **do**

Alg predicts with density matrix \mathbf{W}_t

Nat reveals dyad $\mathbf{x}_t \mathbf{x}_t^\top$

Alg incurs loss $-\mathbf{x}_t^\top \log(\mathbf{W}_t) \mathbf{x}_t$

end for

Density matrix prediction

for trial $t = 1, 2, \dots$ **do**

Alg predicts with density matrix \mathbf{W}_t

Nat reveals dyad $\mathbf{x}_t \mathbf{x}_t^\top$

Alg incurs loss $-\mathbf{x}_t^\top \log(\mathbf{W}_t) \mathbf{x}_t$

end for

for trial $t = 1, 2, \dots$ **do**

Alg predicts with distr. ω_t

Nat reveals $\mathbf{x}_t \in \{e_1, \dots, e_n\}$

Alg incurs loss $-\log(\omega_t^\top \mathbf{x}_t)$

end for

Density matrix prediction

for trial $t = 1, 2, \dots$ **do**

Alg predicts with density matrix W_t

Nat reveals dyad $x_t x_t^\top$

Alg incurs loss $-x_t^\top \log(W_t) x_t$

end for

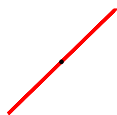
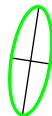
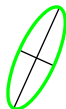
for trial $t = 1, 2, \dots$ **do**

Alg predicts with distr. ω_t

Nat reveals $x_t \in \{e_1, \dots, e_n\}$

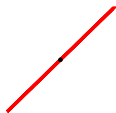
Alg incurs loss $-\log(\omega_t^\top x_t)$

end for



The outcomes: dyads

- A **dyad** $\mathbf{x}\mathbf{x}^\top$ is a rank-one matrix, where \mathbf{x} is a vector in \mathbb{R}^n of unit length.



- A dyad is a **classical outcome in an arbitrary orthonormal basis**:

$$\mathbf{x}\mathbf{x}^\top = \mathbf{U}^\top \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix} \mathbf{U}$$

- There are continuously many dyads.

The predictions: density matrices

- A **density matrix** \mathbf{W} is a convex combination of dyads.
 - \mathbf{W} is a positive-semidefinite matrix of unit trace.
- A density matrix is a **probability vector in an arbitrary orthonormal basis**:

$$\mathbf{W} = \sum_{i=1}^n \omega_i \mathbf{a}_i \mathbf{a}_i^\top$$

eigenvalues ω probability vector
eigenvectors \mathbf{a}_i orthonormal system



The loss: matrix log loss

- The **logarithm** of a density matrix $\mathbf{W} = \sum_i \omega_i \mathbf{a}_i \mathbf{a}_i^\top$ is defined by

$$\log(\mathbf{W}) = \sum_i \log(\omega_i) \mathbf{a}_i \mathbf{a}_i^\top.$$

- Discrepancy between prediction \mathbf{W} and dyad $\mathbf{x}\mathbf{x}^\top$: **matrix log loss**

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

- If **Alg** and **Nat** play in the same eigensystem, i.e. $\mathbf{x} = \mathbf{a}_j$, then **matrix log loss becomes classical log loss**:

$$-\mathbf{x}^\top \log(\mathbf{W}) \mathbf{x} = -\mathbf{a}_j^\top \sum_i \log(\omega_i) \mathbf{a}_i \mathbf{a}_i^\top \mathbf{a}_j = -\log(\omega_j) = -\log(\boldsymbol{\omega}^\top \mathbf{x})$$

Matrix log loss is proper

- The Von Neumann or Quantum entropy:

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

equals the Shannon entropy of eigenvalues α of \mathbf{A} .

- We now compete with the empirical Von Neumann entropy:

$$\inf_{\mathbf{W}} \sum_{t=1}^T -\mathbf{x}_t^\top \log(\mathbf{W}) \mathbf{x}_t = T H(\mathbf{W}^*) \quad \text{where} \quad \mathbf{W}^* = \frac{\sum_{t=1}^T \mathbf{x}_t \mathbf{x}_t^\top}{T}$$

Outline

- 1 Introduction
- 2 Classical Log Loss
- 3 Matrix Log Loss
- 4 Free Matrix Lunch**
- 5 Summary and Open Questions

- Matrix Laplace:

$$\mathbf{W}_{t+1} := \frac{\sum_{q=1}^t \mathbf{x}_q \mathbf{x}_q^\top + \mathbf{I}}{t + n}$$

- Matrix Krychevsky-Trofimoff (KT):

$$\mathbf{W}_{t+1} := \frac{\sum_{q=1}^t \mathbf{x}_q \mathbf{x}_q^\top + \mathbf{I}/2}{t + n/2}$$

- Matrix Laplace:

$$\mathbf{W}_{t+1} := \frac{\sum_{q=1}^t \mathbf{x}_q \mathbf{x}_q^\top + \mathbf{I}}{t + n}$$

$$\boldsymbol{\omega}_{t+1} := \frac{\sum_{q=1}^t \mathbf{x}_q + \mathbf{1}}{t + n}$$

- Matrix Krychevsky-Trofimoff (KT):

$$\mathbf{W}_{t+1} := \frac{\sum_{q=1}^t \mathbf{x}_q \mathbf{x}_q^\top + \mathbf{I}/2}{t + n/2}$$

$$\boldsymbol{\omega}_{t+1} := \frac{\sum_{q=1}^t \mathbf{x}_q + \mathbf{1}/2}{t + n/2}$$

Two Free Matrix Lunches

Theorem

Classical and matrix worst-case regrets coincide for Laplace and for KT.

Two Free Matrix Lunches

Theorem

Classical and matrix worst-case regrets coincide for Laplace and for KT.

But why...?

- If **Alg** plays Laplace or KT, then **Nat** will never go out-eigensystem:
Any sequence of dyads not in same eigensystem is suboptimal for **Nat**
- The classical case is the worst case. No additional regret.
- We learn eigenvectors for free!

Free matrix lunch for Shtarkov?

- Are the classical and matrix prediction games equally hard?
- **Ultimate open problem:** is the *classical minimax regret*

$$\min_{\omega_1} \max_{\mathbf{x}_1} \cdots \min_{\omega_T} \max_{\mathbf{x}_T} \sum_{t=1}^T -\log(\omega_t^\top \mathbf{x}_t) - T H\left(\frac{\sum_{t=1}^T \mathbf{x}_t}{T}\right)$$

equal to the *matrix minimax regret*

$$\min_{\mathbf{W}_1} \max_{\mathbf{x}_1} \cdots \min_{\mathbf{W}_T} \max_{\mathbf{x}_T} \sum_{t=1}^T -\mathbf{x}_t^\top \log(\mathbf{W}_t) \mathbf{x}_t - T H\left(\frac{\sum_{t=1}^T \mathbf{x}_t \mathbf{x}_t^\top}{T}\right)$$

Is there a free matrix lunch for matrix Shtarkov?

- Only numerical evidence for this claim and intermediate conjectures.
- Regret bounds for classical and matrix **Last Step Minimax** coincide.

Outline

- 1 Introduction
- 2 Classical Log Loss
- 3 Matrix Log Loss
- 4 Free Matrix Lunch
- 5 Summary and Open Questions**

Summary

- Matrix extensions of classical algorithms for log loss.
- Learning a matrix of n^2 parameters with regret for n
- Eigenvectors are learned for free
- Classical data is worst-case

Open questions

- Does the free matrix lunch hold for the matrix minimax algorithm?
- A generic method for promoting classical strategies to the matrix domain.
- Different loss functions.

Thank you!