

The Free Matrix Lunch

Wouter M. Koolen



Wojciech Kotłowski



Manfred K. Warmuth



Tuesday 24th April, 2012

The open problem (Warmuth, COLT 2007)

Recent interest in matrix generalizations of classical prediction tasks:

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

The open problem (Warmuth, COLT 2007)

Recent interest in matrix generalizations of classical prediction tasks:

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

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

The open problem (Warmuth, COLT 2007)

Recent interest in matrix generalizations of classical prediction tasks:

- Matrix Hedge (PCA)
- Matrix Winnow (learning subspaces)
- Matrix 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!

The open problem (Warmuth, COLT 2007)

Recent interest in matrix generalizations of classical prediction tasks:

- Matrix Hedge (PCA)
- Matrix Winnow (learning subspaces)
- Matrix 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!

What is going on? Are classical bounds loose, or is there a



Free Matrix Lunch?

This talk I

Fundamental task: **predicting n -ary sequence with logarithmic loss**

- Strong intuition from several interpretations
probability forecasting - data compression - investment
- Algorithms derived from various principles
Bayesian inference - universal coding - optimization - minimax
- Popular. Extremely well-studied. Simple. Often one-line proofs.

This talk I

Fundamental task: **predicting n -ary sequence with logarithmic loss**

- Strong intuition from several interpretations
probability forecasting - data compression - investment
- Algorithms derived from various principles
Bayesian inference - universal coding - optimization - minimax
- Popular. Extremely well-studied. Simple. Often one-line proofs.

We generalise the *problem* and lift the *algorithms* to the matrix domain.

This talk I

Fundamental task: **predicting n -ary sequence with logarithmic loss**

- Strong intuition from several interpretations
probability forecasting - data compression - investment
- Algorithms derived from various principles
Bayesian inference - universal coding - optimization - minimax
- Popular. Extremely well-studied. Simple. Often one-line proofs.

We generalise the *problem* and lift the *algorithms* to the matrix domain.

We prove and explain a



This talk II

We then consider the second fundamental **Hedge** or **dot loss** setting.

This talk II

We then consider the second fundamental **Hedge** or **dot loss** setting.

Here we show matrix prediction is *strictly harder*



Outline

- 1 Introduction
- 2 Classical Log Loss**
- 3 Matrix Log Loss
- 4 Trace Loss Counterexample
- 5 Conclusion

Probability vector prediction

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

Alg predicts with probability vector ω_t

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

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

end for

Probability vector prediction

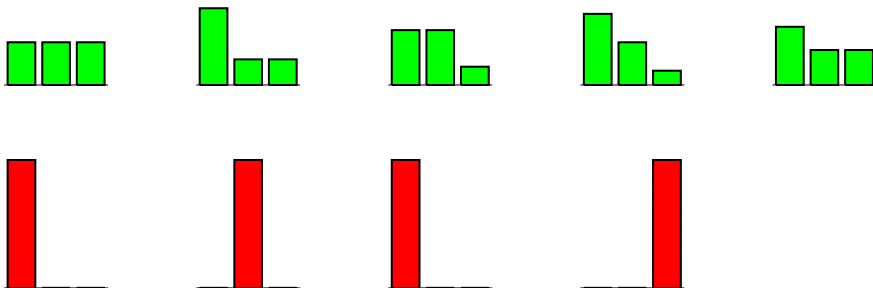
for trial $t = 1, 2, \dots$ **do**

Alg predicts with probability vector ω_t

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

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

end for



Evaluation

Regret is loss of **Alg** minus the loss of the **best fixed prediction**:

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

Evaluation

Regret is loss of **Alg** minus the loss of the **best fixed prediction**:

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

In this problem we compete with the *empirical Shannon entropy*:

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

ω^* is the *maximum likelihood estimator*

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

Algorithms

For the **Laplace** predictor

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

whereas for the **Krychevsky-Trofimoff** predictor

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

Algorithms

For the **Laplace** predictor

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

whereas for the **Krychevsky-Trofimoff** predictor

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

Other algorithms include

$$\textit{Last Step Minimax} \quad \mathcal{R}_T \leq \frac{n-1}{2} \log(T+1) + 1$$

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

Outline

- 1 Introduction
- 2 Classical Log Loss
- 3 Matrix Log Loss**
- 4 Trace Loss Counterexample
- 5 Conclusion

Density matrix prediction

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

Alg predicts with density matrix \mathbf{W}_t

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

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

end for

Density matrix prediction

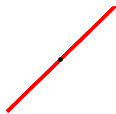
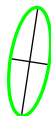
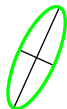
for trial $t = 1, 2, \dots$ **do**

Alg predicts with density matrix \mathbf{W}_t

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

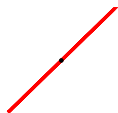
Alg incurs loss $-\mathbf{x}_t^\top \mathbf{log}(\mathbf{W}_t) \mathbf{x}_t$

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.

Positive-semidefinite matrix \mathbf{W} of unit trace

The predictions: density matrices

A **density matrix** \mathbf{W} is a convex combination of dyads.

Positive-semidefinite matrix \mathbf{W} of unit trace

A density matrix is a **probability vector in an arbitrary orthonormal basis**:

Decomposition:

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



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

The predictions: density matrices

A **density matrix** \mathbf{W} is a convex combination of dyads.

Positive-semidefinite matrix \mathbf{W} of unit trace

A density matrix is a **probability vector in an arbitrary orthonormal basis**:

Decomposition:

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



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

Note: different convex combinations of dyads may result in the same density matrix.

The loss: matrix log loss

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

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

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

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

The classical case

If **Alg** and **Nat** play in the same eigensystem, say

$$\mathbf{W} = \sum_i \omega_i \mathbf{e}_i \mathbf{e}_i^\top \quad \text{and} \quad \mathbf{x} = \mathbf{e}_j$$

then matrix log loss becomes classical log loss

$$-\mathbf{x}^\top \mathbf{log}(\mathbf{W}) \mathbf{x} = -\mathbf{x}^\top \sum_i \log(\omega_i) \mathbf{e}_i \mathbf{e}_i^\top \mathbf{x} = -\log(\alpha_j) = -\log(\boldsymbol{\omega}^\top \mathbf{x})$$

The classical case

If **Alg** and **Nat** play in the same eigensystem, say

$$\mathbf{W} = \sum_i \omega_i \mathbf{e}_i \mathbf{e}_i^\top \quad \text{and} \quad \mathbf{x} = \mathbf{e}_j$$

then matrix log loss becomes classical log loss

$$-\mathbf{x}^\top \mathbf{log}(\mathbf{W}) \mathbf{x} = -\mathbf{x}^\top \sum_i \log(\omega_i) \mathbf{e}_i \mathbf{e}_i^\top \mathbf{x} = -\log(\alpha_j) = -\log(\boldsymbol{\omega}^\top \mathbf{x})$$

But both players can deviate. Who is to gain?

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*:

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

Matrix Algorithms

Matrix Laplace predicts with $\mathbf{W}_{t+1} =$

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

and **Matrix KT** predicts with $\mathbf{W}_{t+1} =$

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

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.

Proof for Laplace.

Let \mathbf{W}_t^* denote the best density matrix for the first t outcomes. The regret of matrix Laplace can be bounded as follows:

$$\mathcal{R}_T = \sum_{t=1}^T \ell(\mathbf{W}_t, \mathbf{x}_t \mathbf{x}_t^\top) - \sum_{t=1}^T \ell(\mathbf{W}_T^*, \mathbf{x}_t \mathbf{x}_t^\top) \leq \sum_{t=1}^T \left(\ell(\mathbf{W}_t, \mathbf{x}_t \mathbf{x}_t^\top) - \ell(\mathbf{W}_t^*, \mathbf{x}_t \mathbf{x}_t^\top) \right). \quad (1)$$

Now consider the t^{th} term in the right-hand sum. With $\mathbf{S}_t = \sum_{q=1}^t \mathbf{x}_q \mathbf{x}_q^\top$

$$-\mathbf{x}_t^\top \left(\log \frac{\mathbf{S}_{t-1} + \mathbf{I}}{t-1+n} - \log \frac{\mathbf{S}_t}{t} \right) \mathbf{x}_t = \log \left(\frac{t-1+n}{t} \right) - \mathbf{x}_t^\top (\log(\mathbf{S}_{t-1} + \mathbf{I}) - \log \mathbf{S}_t) \mathbf{x}_t.$$

The matrix part is non-positive since $\mathbf{S}_{t-1} + \mathbf{I} \succeq \mathbf{S}_t$, and the logarithm is matrix monotone. It is zero for any sequence of identical dyads and (1) holds with equality since $\mathbf{W}_t^* = \mathbf{W}_T^*$ for all $t \leq T$. The same upper bound is also met by classical Laplace on any sequence of identical outcomes. □

But why...?

If **Alg** plays Laplace or KT, then **Nat** will never go out-eigensystem.

The classical case is the worst case.

Best eigenvectors are *observed*. No regret.

But why...?

If **Alg** plays Laplace or KT, then **Nat** will never go out-eigensystem.

The classical case is the worst case.

Best eigenvectors are *observed*. No regret.

If **Alg** plays in the eigensystem of past data, will **Nat** do too?

- For matrix log loss: only pathological counterexamples
- For other losses: real counterexamples

Shtarkov: The Queen of Lunches

In the classical case the minimax algorithm is due to Shtarkov.

Ultimate open problem: is the *classical minimax regret*

$$\min_{\omega_1} \max_{x_1} \cdots \min_{\omega_T} \max_{x_T} \sum_{t=1}^T -\log(\omega_t^\top x_t) - T H\left(\frac{\sum_{t=1}^T x_t}{T}\right)$$

equal to the *matrix minimax regret*

$$\min_{W_1} \max_{x_1} \cdots \min_{W_T} \max_{x_T} \sum_{t=1}^T -x_t^\top \mathbf{log}(W_t) x_t - T H\left(\frac{\sum_{t=1}^T x_t x_t^\top}{T}\right)$$

Only numerical evidence for this claim and intermediate conjectures.

Outline

- 1 Introduction
- 2 Classical Log Loss
- 3 Matrix Log Loss
- 4 Trace Loss Counterexample**
- 5 Conclusion

The story for trace loss

The Free Matrix Lunch depends on the choice of loss.

The **dot loss** generalises to the **trace loss**:

$$\ell(\boldsymbol{\omega}, \boldsymbol{l}) = \boldsymbol{\omega}^\top \boldsymbol{l} \qquad \ell(\mathbf{W}, \mathbf{L}) = \text{tr}(\mathbf{W}\mathbf{L})$$

for $\boldsymbol{l} \in \{0, 1\}^n$ and symmetric \mathbf{L} with eigenvalues in $\{0, 1\}^n$.

The story for trace loss

The Free Matrix Lunch depends on the choice of loss.

The **dot loss** generalises to the **trace loss**:

$$\ell(\boldsymbol{\omega}, \boldsymbol{l}) = \boldsymbol{\omega}^\top \boldsymbol{l} \qquad \ell(\mathbf{W}, \mathbf{L}) = \text{tr}(\mathbf{W}\mathbf{L})$$

for $\boldsymbol{l} \in \{0, 1\}^n$ and symmetric \mathbf{L} with eigenvalues in $\{0, 1\}^n$.

Hedge and **Matrix Hedge**

$$\boldsymbol{\omega}_t = \boldsymbol{\omega}_{t-1} e^{-\eta \boldsymbol{l}_t} / Z_t \qquad \mathbf{W}_t = \exp(\log \mathbf{W}_{t-1} - \eta \mathbf{L}_t) / Z_t$$

The story for trace loss

The Free Matrix Lunch depends on the choice of loss.

The **dot loss** generalises to the **trace loss**:

$$\ell(\omega, l) = \omega^\top l \qquad \ell(W, L) = \text{tr}(WL)$$

for $l \in \{0, 1\}^n$ and symmetric L with eigenvalues in $\{0, 1\}^n$.

Hedge and **Matrix Hedge**

$$\omega_t = \omega_{t-1} e^{-\eta l_t} / Z_t \qquad W_t = \exp(\log W_{t-1} - \eta L_t) / Z_t$$

with tuned *learning rate* η both have regrets bounded by

$$\sqrt{\frac{T \log n}{2}} \quad \text{as well as} \quad \sqrt{2L^* \log n} + \log n.$$

No free lunch for trace loss

In dimension $n = 2$ the minimax regrets for T trials are

$$\sqrt{\frac{T+1}{2\pi}} \qquad \sqrt{\frac{T}{4}}$$

Far from a free lunch.

No free lunch for trace loss

In dimension $n = 2$ the minimax regrets for T trials are

$$\sqrt{\frac{T+1}{2\pi}} \qquad \sqrt{\frac{T}{4}}$$

Far from a free lunch.

The matrix case is *different* and *easier* to analyse:

- **Alg** predicts *deterministically* in past eigensystem
- **Nat** plays Hadamard (**worse than classical randomisation**)

No free lunch for trace loss

In dimension $n = 2$ the minimax regrets for T trials are

$$\sqrt{\frac{T+1}{2\pi}} \quad \sqrt{\frac{T}{4}}$$

Far from a free lunch.

The matrix case is *different* and *easier* to analyse:

- **Alg** predicts *deterministically* in past eigensystem
- **Nat** plays Hadamard (**worse than classical randomisation**)

We submitted the case $n > 2$ as an open problem to COLT 2012.

Outline

- 1 Introduction
- 2 Classical Log Loss
- 3 Matrix Log Loss
- 4 Trace Loss Counterexample
- 5 Conclusion**

Summary

- Matrix log loss
 - Learning a matrix of n^2 parameters with regret for n
 - Eigenvectors are learned for free
 - Classical data is worst-case
- Trace loss
 - No free matrix lunch
 - **Nat** exploits matrix power

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 log loss prediction?

Thank you!