

# Minimax Strategies for Square Loss Games

Alan Malek  
University of California, Berkeley

July 28th, 2016

Joint work with:  
Peter L. Bartlett, Wouter M. Koolen, Eiji Takimoto, Manfred  
Warmuth, Yasin Abbasi-Yadkori

## Square loss protocol

Convex set  $\mathcal{C}$ , length  $T$ , and know loss functions  $\ell$ .

For each round  $t = 1, \dots, T$ ,

- ▶ We play  $\mathbf{a}_t \in \mathcal{C}$
- ▶ Nature reveals  $\mathbf{y}_t \in \mathcal{C}$
- ▶ We incur loss

$$\ell(\mathbf{a}_t, \mathbf{y}_t) = \|\mathbf{a}_t - \mathbf{y}_t\|^2$$

For some comparator class  $\mathcal{A}$ , the best comparator is

$$L_T^*(\mathbf{y}_1^T) = \min_{\mathbf{a} \in \mathcal{A}} \sum_{t=1}^T \ell(\mathbf{a}, \mathbf{y}_t).$$

Goal: find a strategy with minimum regret

$$\text{Regret} := \sum_{t=1}^T \ell(\mathbf{a}_t, \mathbf{y}_t) - L_T^*(\mathbf{y}_1^T)$$

# What is minimax?

We play to minimize the worst-case regret. Value is

$$\begin{aligned} V &:= \inf_{\text{Strategies } \mathcal{S}} \sup_{\text{Data } \mathcal{D}} \text{Regret}(\mathcal{S}, \mathcal{D}) \\ &= \text{Regret}(\mathbf{a}_1^T, \mathbf{y}_1^T) \end{aligned}$$

# What is minimax?

We play to minimize the worst-case regret. Value is

$$\begin{aligned} V &:= \inf_{\text{Strategies } \mathcal{S}} \sup_{\text{Data } \mathcal{D}} \text{Regret}(\mathcal{S}, \mathcal{D}) \\ &= \min_{\mathbf{a}_T} \max_{\mathbf{y}_T} \text{Regret}(\mathbf{a}_1^T, \mathbf{y}_1^T) \end{aligned}$$

# What is minimax?

We play to minimize the worst-case regret. Value is

$$\begin{aligned} V &:= \inf_{\text{Strategies } \mathcal{S}} \sup_{\text{Data } \mathcal{D}} \text{Regret}(\mathcal{S}, \mathcal{D}) \\ &= \min_{\mathbf{a}_{T-1} \mathbf{y}_{T-1}} \max_{\mathbf{a}_T \mathbf{y}_T} \min \max \text{Regret}(\mathbf{a}_1^T, \mathbf{y}_1^T) \end{aligned}$$

# What is minimax?

We play to minimize the worst-case regret. Value is

$$\begin{aligned} V &:= \inf_{\text{Strategies } \mathcal{S}} \sup_{\text{Data } \mathcal{D}} \text{Regret}(\mathcal{S}, \mathcal{D}) \\ &= \min_{\mathbf{a}_1} \max_{\mathbf{y}_1} \dots \min_{\mathbf{a}_{T-1}} \max_{\mathbf{y}_{T-1}} \min_{\mathbf{a}_T} \max_{\mathbf{y}_T} \text{Regret}(\mathbf{a}_1^T, \mathbf{y}_1^T) \end{aligned}$$

# What is minimax?

We play to minimize the worst-case regret. Value is

$$\begin{aligned} V &:= \inf_{\text{Strategies } \mathcal{S}} \sup_{\text{Data } \mathcal{D}} \text{Regret}(\mathcal{S}, \mathcal{D}) \\ &= \min_{\mathbf{a}_1} \max_{\mathbf{y}_1} \dots \min_{\mathbf{a}_{T-1}} \max_{\mathbf{y}_{T-1}} \min_{\mathbf{a}_T} \max_{\mathbf{y}_T} \text{Regret}(\mathbf{a}_1^T, \mathbf{y}_1^T) \end{aligned}$$

- ▶ Optimal algorithm against worst case adversary
- ▶ How can we compute this?
- ▶ Backwards induction / dynamic programming

## Value-to-go

Consider what happens after  $t$  rounds:

$$\begin{aligned} V &= \min_{\mathbf{a}_1} \max_{\mathbf{y}_1} \dots \min_{\mathbf{a}_T} \max_{\mathbf{y}_T} \sum_{t=1}^T \ell(\mathbf{a}_t, \mathbf{y}_t) - L_T^*(\mathbf{y}_1^T) \\ &= \min_{\mathbf{a}_1} \max_{\mathbf{y}_1} \dots \min_{\mathbf{a}_t} \max_{\mathbf{y}_t} \sum_{\tau=1}^t \ell(\mathbf{a}_\tau, \mathbf{y}_\tau) \\ &\quad + \underbrace{\min_{\mathbf{a}_{t+1}} \max_{\mathbf{y}_{t+1}} \dots \min_{\mathbf{a}_T} \max_{\mathbf{y}_T} \sum_{\tau=t+1}^T \ell(\mathbf{a}_\tau, \mathbf{y}_\tau) - L_T^*(\mathbf{y}_1^T)}_{:= V_t(\mathbf{y}_1^t), \text{ the value-to-go with history } \mathbf{a}_1^t, \mathbf{y}_1^t} \end{aligned}$$



## Value-to-go

Consider what happens after  $t$  rounds:

$$\begin{aligned} V &= \min_{\mathbf{a}_1} \max_{\mathbf{y}_1} \dots \min_{\mathbf{a}_T} \max_{\mathbf{y}_T} \sum_{t=1}^T \ell(\mathbf{a}_t, \mathbf{y}_t) - L_T^*(\mathbf{y}_1^T) \\ &= \min_{\mathbf{a}_1} \max_{\mathbf{y}_1} \dots \min_{\mathbf{a}_t} \max_{\mathbf{y}_t} \sum_{\tau=1}^t \ell(\mathbf{a}_\tau, \mathbf{y}_\tau) \\ &\quad + \underbrace{\min_{\mathbf{a}_{t+1}} \max_{\mathbf{y}_{t+1}} \dots \min_{\mathbf{a}_T} \max_{\mathbf{y}_T} \sum_{\tau=t+1}^T \ell(\mathbf{a}_\tau, \mathbf{y}_\tau) - L_T^*(\mathbf{y}_1^T)}_{:= V_t(\mathbf{y}_1^t), \text{ the value-to-go with history } \mathbf{a}_1^t, \mathbf{y}_1^t} \end{aligned}$$

Inductive definition:

$$V_T(\mathbf{y}_1^T) := -L_T^*(\mathbf{y}_1^T) \quad (1)$$

$$V_{t-1}(\mathbf{y}_1^{t-1}) := \min_{\mathbf{a}_t} \max_{\mathbf{y}_t} \ell(\mathbf{a}_t, \mathbf{y}_t) + V_t(\mathbf{y}_1, \dots, \mathbf{y}_t) \quad (2)$$

# Value-to-go

The minimax regret  $V$  equals value-to-go  $V_0(\epsilon)$  (empty history).

The minimax strategy: after seeing  $y_1, \dots, y_{t-1}$ ,

- ▶ Compute  $V_t(y_1, \dots, y_t)$
- ▶ Choose  $a_t$  as the minimizer of

$$V(y_1, \dots, y_{t-1}) := \min_{a_t} \max_{y_t} \ell(a_t, y_t) + V(y_1, \dots, y_t)$$

# Value-to-go

The minimax regret  $V$  equals value-to-go  $V_0(\epsilon)$  (empty history).

The minimax strategy: after seeing  $y_1, \dots, y_{t-1}$ ,

- ▶ Compute  $V_t(y_1, \dots, y_t)$
- ▶ Choose  $a_t$  as the minimizer of

$$V(y_1, \dots, y_{t-1}) := \min_{a_t} \max_{y_t} \ell(a_t, y_t) + V(y_1, \dots, y_t)$$

Problem: this is expensive (usually exponentially so).

# Outline

- ▶ What is minimax?
- ▶ Two minimax square loss games
- ▶ Mimimax fixed-design online linear regression
- ▶ Minimax time series prediction

## Section 1

### Square loss game

## Square loss protocol (with Koolen and Bartlett)

Convex set  $\mathcal{C}$ , length  $T$ , and know loss functions  $\ell$ .

For each round  $t = 1, \dots, T$ ,

- ▶ We play  $\mathbf{a}_t \in \mathcal{C}$
- ▶ Nature reveals  $\mathbf{y}_t \in \mathcal{C}$
- ▶ We incur loss

matrix  $\mathbf{W}$  weights prediction errors



$$\ell(\mathbf{a}_t, \mathbf{y}_t) := \|\mathbf{a}_t - \mathbf{y}_t\|_{\mathbf{W}}^2 = (\mathbf{a}_t - \mathbf{y}_t)^\top \mathbf{W}^{-1} (\mathbf{a}_t - \mathbf{y}_t)$$

Our goal is to minimize regret w.r.t. best fixed action  $\mathbf{a}$  in hindsight

$$\text{Regret} := \sum_{t=1}^T \ell(\mathbf{a}_t, \mathbf{y}_t) - \min_{\mathbf{a}} \sum_{t=1}^T \ell(\mathbf{a}, \mathbf{y}_t)$$

Notation:  $\mathbf{a}_1^t = (\mathbf{a}_1, \dots, \mathbf{a}_t)$ .

## Solving the minimax strategy

- ▶ Using sufficient statistics

$$\mathbf{s}_t = \sum_{\tau=1}^t \mathbf{y}_\tau \quad \text{and} \quad \sigma_t^2 = \sum_{\tau=1}^t \mathbf{y}_\tau^\top \mathbf{W}^{-1} \mathbf{y}_\tau$$

- ▶ First, we need  $L_T^*(\mathbf{y}_1^T)$ :

$$L_T^* = \inf_{\mathbf{a} \in \mathbb{R}^d} \sum_{t=1}^T \|\mathbf{a} - \mathbf{y}_t\|_{\mathbf{W}}^2 = \sigma_T^2 - \frac{1}{T} \mathbf{s}_T^\top \mathbf{W}^{-1} \mathbf{s}_T$$

and the minimizer is the mean outcome  $\mathbf{a}^* = \frac{1}{T} \sum_{t=1}^T \mathbf{y}_t$ .

## Calculating the value function for $\mathcal{C} = \triangle$

- ▶ Need to solve the backwards induction
- ▶ Base case:  $V_T(y_1^T) = -L_T^* = \frac{1}{T} s_T^T W^{-1} s_T - \sigma^2_T$



## Calculating the value function for $\mathcal{C} = \triangle$

- ▶ Need to solve the backwards induction
- ▶ Base case:  $V_T(y_1^T) = -L_T^* = \frac{1}{T} \mathbf{s}_T^T \mathbf{W}^{-1} \mathbf{s}_T - \sigma_T^2$
- ▶ “Guess”:

$$V_t(\mathbf{s}_t, \sigma_t^2) = \alpha_t \mathbf{s}_t^T \mathbf{W}^{-1} \mathbf{s}_t - \sigma_t^2 + (1 - t\alpha_t) \text{diag}(\mathbf{W}^{-1})^T \mathbf{s}_t + \gamma_t,$$

## Calculating the value function for $\mathcal{C} = \triangle$

- ▶ Need to solve the backwards induction
- ▶ Base case:  $V_T(y_1^T) = -L_T^* = \frac{1}{T} \mathbf{s}_T^T \mathbf{W}^{-1} \mathbf{s}_T - \sigma^2_T$
- ▶ “Guess”:

$$V_t(\mathbf{s}_t, \sigma^2_t) = \alpha_t \mathbf{s}_t^T \mathbf{W}^{-1} \mathbf{s}_t - \sigma^2_t + (1 - t\alpha_t) \text{diag}(\mathbf{W}^{-1})^T \mathbf{s}_t + \gamma_t,$$

- ▶ Base case:  $\alpha_T = \frac{1}{T}$ ,  $\gamma_t = 0$

## Calculating the value function for $\mathcal{C} = \triangle$

- ▶ Need to solve the backwards induction
- ▶ Base case:  $V_T(y_1^T) = -L_T^* = \frac{1}{T} \mathbf{s}_T^T \mathbf{W}^{-1} \mathbf{s}_T - \sigma^2_T$
- ▶ “Guess”:

$$V_t(\mathbf{s}_t, \sigma^2_t) = \alpha_t \mathbf{s}_t^T \mathbf{W}^{-1} \mathbf{s}_t - \sigma^2_t + (1 - t\alpha_t) \text{diag}(\mathbf{W}^{-1})^T \mathbf{s}_t + \gamma_t,$$

- ▶ Base case:  $\alpha_T = \frac{1}{T}$ ,  $\gamma_t = 0$
- ▶ Induction:

$$V_t(\mathbf{s}_t, \sigma^2_t) = \min_{\mathbf{a} \in \Delta} \max_{\mathbf{y} \in \Delta} \ell(\mathbf{a}, \mathbf{y}) + V_{t+1}(\mathbf{s}_t + \mathbf{y}, \sigma^2_t + \mathbf{y}^T \mathbf{W}^{-1} \mathbf{y})$$

$$\begin{aligned}
V_t(\mathbf{s}_t, \sigma_t^2) &= \min_{\mathbf{a} \in \Delta} \max_{\mathbf{y} \in \Delta} \|\mathbf{a} - \mathbf{y}\|_{\mathbf{W}}^2 + \alpha_t (\mathbf{s}_t + \mathbf{y})^\top \mathbf{W}^{-1} (\mathbf{s}_t + \mathbf{y}) \\
&\quad - (\sigma_t^2 + \mathbf{y}^\top \mathbf{W}^{-1} \mathbf{y}) + \gamma_t \\
&\quad + (1 - t\alpha_t) \text{diag}(\mathbf{W}^{-1})^\top (\mathbf{s}_t + \mathbf{y})
\end{aligned}$$

$$\begin{aligned}
V_t(\mathbf{s}_t, \sigma_t^2) &= \min_{\mathbf{a} \in \Delta} \max_{\mathbf{y} \in \Delta} \|\mathbf{a} - \mathbf{y}\|_{\mathbf{W}}^2 + \alpha_t (\mathbf{s}_t + \mathbf{y})^\top \mathbf{W}^{-1} (\mathbf{s}_t + \mathbf{y}) \\
&\quad - (\sigma_t^2 + \mathbf{y}^\top \mathbf{W}^{-1} \mathbf{y}) + \gamma_t \\
&\quad + (1 - t\alpha_t) \text{diag}(\mathbf{W}^{-1})^\top (\mathbf{s}_t + \mathbf{y}) \\
&= \min_{\mathbf{a} \in \Delta} \max_{\mathbf{y} \in \Delta} \|\mathbf{a} - \mathbf{y}\|_{\mathbf{W}}^2 + (\alpha_t - 1) \mathbf{y}^\top \mathbf{W}^{-1} \mathbf{y} \\
&\quad + \underbrace{(2\alpha_t \mathbf{W}^{-1} \mathbf{s}_t + (1 - t\alpha_t) \text{diag}(\mathbf{W}^{-1}))^\top}_{:= \mathbf{b}^\top} \mathbf{y} + c
\end{aligned}$$

$$\begin{aligned}
V_t(\mathbf{s}_t, \sigma_t^2) &= \min_{\mathbf{a} \in \Delta} \max_{\mathbf{y} \in \Delta} \|\mathbf{a} - \mathbf{y}\|_{\mathbf{W}}^2 + \alpha_t (\mathbf{s}_t + \mathbf{y})^\top \mathbf{W}^{-1} (\mathbf{s}_t + \mathbf{y}) \\
&\quad - (\sigma_t^2 + \mathbf{y}^\top \mathbf{W}^{-1} \mathbf{y}) + \gamma_t \\
&\quad + (1 - t\alpha_t) \text{diag}(\mathbf{W}^{-1})^\top (\mathbf{s}_t + \mathbf{y}) \\
&= \min_{\mathbf{a} \in \Delta} \max_{\mathbf{y} \in \Delta} \|\mathbf{a} - \mathbf{y}\|_{\mathbf{W}}^2 + (\alpha_t - 1) \mathbf{y}^\top \mathbf{W}^{-1} \mathbf{y} \\
&\quad + \underbrace{(2\alpha_t \mathbf{W}^{-1} \mathbf{s}_t + (1 - t\alpha_t) \text{diag}(\mathbf{W}^{-1}))^\top}_{:= \mathbf{b}^\top} \mathbf{y} + c \\
&= \min_{\mathbf{a} \in \Delta} \max_{\mathbf{y} \in \Delta} \|\mathbf{a} - \mathbf{y}\|_{\mathbf{W}}^2 + (\alpha_t - 1) \mathbf{y}^\top \mathbf{W}^{-1} \mathbf{y} + \mathbf{b}^\top \mathbf{y} + c
\end{aligned}$$

$$\begin{aligned}
V_t(\mathbf{s}_t, \sigma_t^2) &= \min_{\mathbf{a} \in \Delta} \max_{\mathbf{y} \in \Delta} \|\mathbf{a} - \mathbf{y}\|_{\mathbf{W}}^2 + \alpha_t (\mathbf{s}_t + \mathbf{y})^\top \mathbf{W}^{-1} (\mathbf{s}_t + \mathbf{y}) \\
&\quad - (\sigma_t^2 + \mathbf{y}^\top \mathbf{W}^{-1} \mathbf{y}) + \gamma_t \\
&\quad + (1 - t\alpha_t) \text{diag}(\mathbf{W}^{-1})^\top (\mathbf{s}_t + \mathbf{y}) \\
&= \min_{\mathbf{a} \in \Delta} \max_{\mathbf{y} \in \Delta} \|\mathbf{a} - \mathbf{y}\|_{\mathbf{W}}^2 + (\alpha_t - 1) \mathbf{y}^\top \mathbf{W}^{-1} \mathbf{y} \\
&\quad + \underbrace{(2\alpha_t \mathbf{W}^{-1} \mathbf{s}_t + (1 - t\alpha_t) \text{diag}(\mathbf{W}^{-1}))^\top}_{:= \mathbf{b}^\top} \mathbf{y} + c \\
&= \min_{\mathbf{a} \in \Delta} \max_{\mathbf{y} \in \Delta} \|\mathbf{a} - \mathbf{y}\|_{\mathbf{W}}^2 + (\alpha_t - 1) \mathbf{y}^\top \mathbf{W}^{-1} \mathbf{y} + \mathbf{b}^\top \mathbf{y} + c \\
&= \min_{\mathbf{a} \in \Delta} \max_k \|\mathbf{a} - \mathbf{e}_k\|_{\mathbf{W}}^2 + (\alpha_t - 1) \mathbf{e}_k^\top \mathbf{W}^{-1} \mathbf{e}_k + \mathbf{b}^\top \mathbf{e}_k + c
\end{aligned}$$

$$\begin{aligned}
V_t(\mathbf{s}_t, \sigma_t^2) &= \min_{\mathbf{a} \in \Delta} \max_{\mathbf{y} \in \Delta} \|\mathbf{a} - \mathbf{y}\|_{\mathbf{W}}^2 + \alpha_t (\mathbf{s}_t + \mathbf{y})^\top \mathbf{W}^{-1} (\mathbf{s}_t + \mathbf{y}) \\
&\quad - (\sigma_t^2 + \mathbf{y}^\top \mathbf{W}^{-1} \mathbf{y}) + \gamma_t \\
&\quad + (1 - t\alpha_t) \text{diag}(\mathbf{W}^{-1})^\top (\mathbf{s}_t + \mathbf{y}) \\
&= \min_{\mathbf{a} \in \Delta} \max_{\mathbf{y} \in \Delta} \|\mathbf{a} - \mathbf{y}\|_{\mathbf{W}}^2 + (\alpha_t - 1) \mathbf{y}^\top \mathbf{W}^{-1} \mathbf{y} \\
&\quad + \underbrace{(2\alpha_t \mathbf{W}^{-1} \mathbf{s}_t + (1 - t\alpha_t) \text{diag}(\mathbf{W}^{-1}))^\top}_{:= \mathbf{b}^\top} \mathbf{y} + c \\
&= \min_{\mathbf{a} \in \Delta} \max_{\mathbf{y} \in \Delta} \|\mathbf{a} - \mathbf{y}\|_{\mathbf{W}}^2 + (\alpha_t - 1) \mathbf{y}^\top \mathbf{W}^{-1} \mathbf{y} + \mathbf{b}^\top \mathbf{y} + c \\
&= \min_{\mathbf{a} \in \Delta} \max_k \|\mathbf{a} - \mathbf{e}_k\|_{\mathbf{W}}^2 + (\alpha_t - 1) \mathbf{e}_k^\top \mathbf{W}^{-1} \mathbf{e}_k + \mathbf{b}^\top \mathbf{e}_k + c \\
&= \max_p \min_{\mathbf{a} \in \Delta} \mathbb{E}_{k \sim p} \left[ \|\mathbf{a} - \mathbf{e}_k\|_{\mathbf{W}}^2 + (\alpha_t - 1) \mathbf{e}_k^\top \mathbf{W}^{-1} \mathbf{e}_k + \mathbf{b}^\top \mathbf{e}_k \right] +
\end{aligned}$$



$$\begin{aligned}
V_t(\mathbf{s}_t, \sigma_t^2) &= \min_{\mathbf{a} \in \Delta} \max_{\mathbf{y} \in \Delta} \|\mathbf{a} - \mathbf{y}\|_{\mathbf{W}}^2 + \alpha_t (\mathbf{s}_t + \mathbf{y})^\top \mathbf{W}^{-1} (\mathbf{s}_t + \mathbf{y}) \\
&\quad - (\sigma_t^2 + \mathbf{y}^\top \mathbf{W}^{-1} \mathbf{y}) + \gamma_t \\
&\quad + (1 - t\alpha_t) \text{diag}(\mathbf{W}^{-1})^\top (\mathbf{s}_t + \mathbf{y}) \\
&= \min_{\mathbf{a} \in \Delta} \max_{\mathbf{y} \in \Delta} \|\mathbf{a} - \mathbf{y}\|_{\mathbf{W}}^2 + (\alpha_t - 1) \mathbf{y}^\top \mathbf{W}^{-1} \mathbf{y} \\
&\quad + \underbrace{(2\alpha_t \mathbf{W}^{-1} \mathbf{s}_t + (1 - t\alpha_t) \text{diag}(\mathbf{W}^{-1}))^\top}_{:= \mathbf{b}^\top} \mathbf{y} + c \\
&= \min_{\mathbf{a} \in \Delta} \max_{\mathbf{y} \in \Delta} \|\mathbf{a} - \mathbf{y}\|_{\mathbf{W}}^2 + (\alpha_t - 1) \mathbf{y}^\top \mathbf{W}^{-1} \mathbf{y} + \mathbf{b}^\top \mathbf{y} + c \\
&= \min_{\mathbf{a} \in \Delta} \max_k \|\mathbf{a} - \mathbf{e}_k\|_{\mathbf{W}}^2 + (\alpha_t - 1) \mathbf{e}_k^\top \mathbf{W}^{-1} \mathbf{e}_k + \mathbf{b}^\top \mathbf{e}_k + c \\
&= \max_{\mathbf{p}} \min_{\mathbf{a} \in \Delta} \mathbb{E}_{k \sim \mathbf{p}} \left[ \|\mathbf{a} - \mathbf{e}_k\|_{\mathbf{W}}^2 + (\alpha_t - 1) \mathbf{e}_k^\top \mathbf{W}^{-1} \mathbf{e}_k + \mathbf{b}^\top \mathbf{e}_k \right] + c \\
&= \max_{\mathbf{p}} -\mathbf{p}^\top \mathbf{W}^{-1} \mathbf{p} + (\alpha_t \text{diag}(\mathbf{W}^{-1}) + \mathbf{b})^\top \mathbf{p} + c
\end{aligned}$$

$$\begin{aligned}
V_t(\mathbf{s}_t, \sigma_t^2) &= \min_{\mathbf{a} \in \Delta} \max_{\mathbf{y} \in \Delta} \|\mathbf{a} - \mathbf{y}\|_{\mathbf{W}}^2 + \alpha_t (\mathbf{s}_t + \mathbf{y})^\top \mathbf{W}^{-1} (\mathbf{s}_t + \mathbf{y}) \\
&\quad - (\sigma_t^2 + \mathbf{y}^\top \mathbf{W}^{-1} \mathbf{y}) + \gamma_t \\
&\quad + (1 - t\alpha_t) \text{diag}(\mathbf{W}^{-1})^\top (\mathbf{s}_t + \mathbf{y}) \\
&= \min_{\mathbf{a} \in \Delta} \max_{\mathbf{y} \in \Delta} \|\mathbf{a} - \mathbf{y}\|_{\mathbf{W}}^2 + (\alpha_t - 1) \mathbf{y}^\top \mathbf{W}^{-1} \mathbf{y} \\
&\quad + \underbrace{(2\alpha_t \mathbf{W}^{-1} \mathbf{s}_t + (1 - t\alpha_t) \text{diag}(\mathbf{W}^{-1}))^\top}_{:= \mathbf{b}^\top} \mathbf{y} + c \\
&= \min_{\mathbf{a} \in \Delta} \max_{\mathbf{y} \in \Delta} \|\mathbf{a} - \mathbf{y}\|_{\mathbf{W}}^2 + (\alpha_t - 1) \mathbf{y}^\top \mathbf{W}^{-1} \mathbf{y} + \mathbf{b}^\top \mathbf{y} + c \\
&= \min_{\mathbf{a} \in \Delta} \max_k \|\mathbf{a} - \mathbf{e}_k\|_{\mathbf{W}}^2 + (\alpha_t - 1) \mathbf{e}_k^\top \mathbf{W}^{-1} \mathbf{e}_k + \mathbf{b}^\top \mathbf{e}_k + c \\
&= \max_{\mathbf{p}} \min_{\mathbf{a} \in \Delta} \mathbb{E}_{k \sim \mathbf{p}} \left[ \|\mathbf{a} - \mathbf{e}_k\|_{\mathbf{W}}^2 + (\alpha_t - 1) \mathbf{e}_k^\top \mathbf{W}^{-1} \mathbf{e}_k + \mathbf{b}^\top \mathbf{e}_k \right] + c \\
&= \max_{\mathbf{p}} -\mathbf{p}^\top \mathbf{W}^{-1} \mathbf{p} + (\alpha_t \text{diag}(\mathbf{W}^{-1}) + \mathbf{b})^\top \mathbf{p} + c
\end{aligned}$$

Easy to solve via Lagrange multipliers.

## Simplex game (e.g. Brier game)

### Theorem

Let  $\mathcal{C} = \Delta$ . For  $\mathbf{W}$  satisfying an alignment condition, the value-to-go is

$$V_t(\mathbf{s}_t, \sigma_t^2) = \alpha_t \mathbf{s}_t^\top \mathbf{W}^{-1} \mathbf{s}_t - \sigma_t^2 + (1 - t\alpha_t) \text{diag}(\mathbf{W}^{-1})^\top \mathbf{s}_t + \text{const}$$

with coefficients

$$\alpha_T = \frac{1}{T} \text{ and } \alpha_t = \alpha_{t+1}^2 + \alpha_{t+1}.$$

The minimax and maximin strategies are

$$\mathbf{a}_t = \mathbf{p}_t = \frac{\mathbf{s}_t}{t} t\alpha_{t+1} + \mathbf{c}(1 - t\alpha_{t+1})$$

which is data mean  $\frac{\mathbf{s}_t}{t}$  shrunk towards center

$$\mathbf{c} = \frac{\mathbf{W}\mathbf{1}}{\mathbf{1}^\top \mathbf{W}\mathbf{1}} + \left( \mathbf{W} - \frac{\mathbf{W}\mathbf{1}\mathbf{1}^\top \mathbf{W}}{\mathbf{1}^\top \mathbf{W}\mathbf{1}} \right) \text{diag}(\mathbf{W}^{-1})$$

# Ball game

## Theorem

Let  $\mathcal{C} = \bigcirc$ . For any positive definite  $\mathbf{W}$  the value-to-go is

$$V_t(\mathbf{s}_t, \sigma_t^2) = \mathbf{s}_t^\top \mathbf{A}_t \mathbf{s} - \sigma_t^2 + \text{const.}$$

For round  $t + 1$ , the minimax strategy plays

$$\mathbf{a}^* = \left( \lambda_{\max} \mathbf{I} - (\mathbf{A}_{t+1} - \mathbf{W}^{-1}) \right)^{-1} \mathbf{A}_{t+1} \mathbf{s}$$

with coefficients  $\mathbf{A}_T = \frac{1}{T} \mathbf{W}^{-1}$  and

$$\mathbf{A}_t = \mathbf{A}_{t+1} \left( \mathbf{W}^{-1} + \lambda_{\max} \mathbf{I} - \mathbf{A}_{t+1} \right)^{-1} \mathbf{A}_{t+1} + \mathbf{A}_{t+1}.$$

# Regret bounds

- ▶  $\text{Regret}_{\text{Brier}} \propto \sum_{t=1}^T \alpha_t.$
- ▶  $\text{Regret}_{\text{Ball}} = \lambda_{\max}(\mathbf{W}^{-1}) \sum_{t=1}^T \alpha_t.$
- ▶ [1] show that  $\sum_{t=1}^T \alpha_t = O(\log(T) - \log \log(T)).$
- ▶ Compare with  $O(\log(T))$  of Follow the Leader.



E. Takimoto, M. Warmuth

The minimax strategy for Gaussian density estimation  
In *COLT '00*

## Section 2

### Online Linear regression

# Online linear regression (with Bartlett, Koolen, Takimoto, Warmuth)

Fix a covariate sequence  $\mathbf{x}_1, \dots, \mathbf{x}_T$  (fixed design) and length  $T$ .  
For each round  $t = 1, \dots, T$ ,

- ▶ We play  $\mathbf{a}_t \in \mathbb{R}$
- ▶ Nature reveals  $y_t \in [-B_t, B_t]$
- ▶ We incur loss

$$\ell(\mathbf{a}_t, y_t) = (\mathbf{a}_t - y_t)^2$$

- ▶ Minimax Regret is

$$\min_{\mathbf{a}_1} \max_{y_1} \cdots \min_{\mathbf{a}_T} \max_{y_T} \sum_{t=1}^T (\mathbf{a}_t - y_t)^2 - \min_{\theta \in \mathbb{R}^d} \sum_{t=1}^T (\theta^\top \mathbf{x}_t - y_t)^2$$

# Online linear regression (with Bartlett, Koolen, Takimoto, Warmuth)

Fix a covariate sequence  $\mathbf{x}_1, \dots, \mathbf{x}_T$  (fixed design) and length  $T$ .  
For each round  $t = 1, \dots, T$ ,

- ▶ We play  $\mathbf{a}_t \in \mathbb{R}$
- ▶ Nature reveals  $y_t \in [-B_t, B_t]$
- ▶ We incur loss

$$\ell(\mathbf{a}_t, y_t) = (\mathbf{a}_t - y_t)^2$$

- ▶ Minimax Regret is

$$\min_{\mathbf{a}_1} \max_{y_1} \cdots \min_{\mathbf{a}_T} \max_{y_T} \underbrace{\sum_{t=1}^T (\mathbf{a}_t - y_t)^2}_{\text{algorithm}} - \min_{\theta \in \mathbb{R}^d} \sum_{t=1}^T (\theta^\top \mathbf{x}_t - y_t)^2$$



# Online linear regression (with Bartlett, Koolen, Takimoto, Warmuth)

Fix a covariate sequence  $\mathbf{x}_1, \dots, \mathbf{x}_T$  (fixed design) and length  $T$ .  
For each round  $t = 1, \dots, T$ ,

- ▶ We play  $\mathbf{a}_t \in \mathbb{R}$
- ▶ Nature reveals  $y_t \in [-B_t, B_t]$
- ▶ We incur loss

$$\ell(\mathbf{a}_t, y_t) = (\mathbf{a}_t - y_t)^2$$

- ▶ Minimax Regret is

$$\min_{\mathbf{a}_1} \max_{y_1} \cdots \min_{\mathbf{a}_T} \max_{y_T} \underbrace{\sum_{t=1}^T (\mathbf{a}_t - y_t)^2}_{\text{algorithm}} - \underbrace{\min_{\theta \in \mathbb{R}^d} \sum_{t=1}^T (\theta^\top \mathbf{x}_t - y_t)^2}_{\text{best linear predictor}}$$

# Offline problem

- ▶ Define

$$\mathbf{s}_t = \sum_{\tau=1}^t y_{\tau} \mathbf{x}_{\tau}, \quad \sigma_t^2 = \sum_{\tau=1}^t y_{\tau}^2, \quad \mathbf{P}_T = \left( \sum_{t=1}^T \mathbf{x}_t \mathbf{x}_t^{\top} \right)^{-1}$$

- ▶ What is the best *linear predictor* in hindsight:

$$\min_{\theta \in \mathbb{R}^d} \sum_{t=1}^T (\theta^{\top} \mathbf{x}_t - y_t)^2?$$

# Offline problem

- ▶ Define

$$\mathbf{s}_t = \sum_{\tau=1}^t y_{\tau} \mathbf{x}_{\tau}, \quad \sigma^2_t = \sum_{\tau=1}^t y_{\tau}^2, \quad \mathbf{P}_T = \left( \sum_{t=1}^T \mathbf{x}_t \mathbf{x}_t^{\top} \right)^{-1}$$

- ▶ What is the best *linear predictor* in hindsight:

$$\min_{\theta \in \mathbb{R}^d} \sum_{t=1}^T (\theta^{\top} \mathbf{x}_t - y_t)^2?$$

- ▶ Ordinary least squares:

$$\theta^* = \mathbf{P}_T \mathbf{s}_T$$

with loss

$$L_T^* = \sigma^2_T - \mathbf{s}_T^{\top} \mathbf{P}_T \mathbf{s}_T.$$

## Various algorithms

Popular approaches:

$$\hat{y}_{t+1}^{\text{FTL}} := \mathbf{x}_{t+1}^\top \left( \sum_{q=1}^t \mathbf{x}_q \mathbf{x}_q^\top \right)^{-1} \mathbf{s}_t$$

$$\hat{y}_{t+1}^{\text{Ridge}} := \mathbf{x}_{t+1}^\top \left( \sum_{q=1}^t \mathbf{x}_q \mathbf{x}_q^\top + \lambda \mathbf{I} \right)^{-1} \mathbf{s}_t$$

$$\hat{y}_{t+1}^{\text{LSM}} := \mathbf{x}_{t+1}^\top \left( \sum_{q=1}^{t+1} \mathbf{x}_q \mathbf{x}_q^\top \right)^{-1} \mathbf{s}_t$$

## Various algorithms

Popular approaches:

$$\hat{y}_{t+1}^{\text{FTL}} := \mathbf{x}_{t+1}^\top \left( \sum_{q=1}^t \mathbf{x}_q \mathbf{x}_q^\top \right)^{-1} \mathbf{s}_t$$

$$\hat{y}_{t+1}^{\text{Ridge}} := \mathbf{x}_{t+1}^\top \left( \sum_{q=1}^t \mathbf{x}_q \mathbf{x}_q^\top + \lambda \mathbf{I} \right)^{-1} \mathbf{s}_t$$

$$\hat{y}_{t+1}^{\text{LSM}} := \mathbf{x}_{t+1}^\top \left( \sum_{q=1}^{t+1} \mathbf{x}_q \mathbf{x}_q^\top \right)^{-1} \mathbf{s}_t$$

Claim:

$$\hat{y}_{t+1}^{\text{MM}} := \mathbf{x}_{t+1}^\top \mathbf{P}_{t+1} \mathbf{s}_t$$

## Value-to-go stays quadratic

- ▶ We show by induction that

$$V_t(\mathbf{s}_t, \sigma_t^2) = \mathbf{s}_t^\top \mathbf{P}_t \mathbf{s}_t - \sigma_t^2 + \gamma_t,$$

with the  $\gamma_t$  coefficients recursively defined by

$$\gamma_T = 0, \quad \gamma_t = \gamma_{t+1} + B_{t+1}^2 \mathbf{x}_{t+1}^\top \mathbf{P}_{t+1} \mathbf{x}_{t+1}$$

## Value-to-go stays quadratic

- ▶ We show by induction that

$$V_t(\mathbf{s}_t, \sigma_t^2) = \mathbf{s}_t^\top \mathbf{P}_t \mathbf{s}_t - \sigma_t^2 + \gamma_t,$$

with the  $\gamma_t$  coefficients recursively defined by

$$\gamma_T = 0, \quad \gamma_t = \gamma_{t+1} + B_{t+1}^2 \mathbf{x}_{t+1}^\top \mathbf{P}_{t+1} \mathbf{x}_{t+1}$$

- ▶ Base case is easy:

$$V_T = -L_T^* = \mathbf{s}_T^\top \mathbf{P}_T \mathbf{s}_T - \sigma_T^2$$

- ▶ Backwards induction gives

$$V_t(\mathbf{s}_t, \sigma_t^2) := \min_{\hat{y}_{t+1}} \max_{y_{t+1}} (\hat{y}_{t+1} - y_{t+1})^2 \\ + V_{t+1}(\mathbf{s}_t + y_{t+1} \mathbf{x}_{t+1}, \sigma_t^2 + y_{t+1}^2),$$



► Backwards induction gives

$$\begin{aligned} V_t(\mathbf{s}_t, \sigma_t^2) &:= \min_{\hat{y}_{t+1}} \max_{y_{t+1}} (\hat{y}_{t+1} - y_{t+1})^2 \\ &\quad + V_{t+1}(\mathbf{s}_t + y_{t+1} \mathbf{x}_{t+1}, \sigma_t^2 + y_{t+1}^2), \\ &= \min_{\hat{y}_{t+1}} \max_{y_{t+1}} (\hat{y}_{t+1} - y_{t+1})^2 \\ &\quad + (\mathbf{s}_t + y_{t+1} \mathbf{x}_{t+1})^\top \mathbf{P}_{t+1} (\mathbf{s}_t + y_{t+1} \mathbf{x}_{t+1}) \\ &\quad - (\sigma_t^2 + y_{t+1}^2) + \gamma_{t+1} \end{aligned}$$

- ▶ Backwards induction gives

$$\begin{aligned}
 V_t(\mathbf{s}_t, \sigma_t^2) &:= \min_{\hat{y}_{t+1}} \max_{y_{t+1}} (\hat{y}_{t+1} - y_{t+1})^2 \\
 &\quad + V_{t+1}(\mathbf{s}_t + y_{t+1} \mathbf{x}_{t+1}, \sigma_t^2 + y_{t+1}^2), \\
 &= \min_{\hat{y}_{t+1}} \max_{y_{t+1}} (\hat{y}_{t+1} - y_{t+1})^2 \\
 &\quad + (\mathbf{s}_t + y_{t+1} \mathbf{x}_{t+1})^\top \mathbf{P}_{t+1} (\mathbf{s}_t + y_{t+1} \mathbf{x}_{t+1}) \\
 &\quad - (\sigma_t^2 + y_{t+1}^2) + \gamma_{t+1}
 \end{aligned}$$

- ▶ This is convex in  $y_{t+1}$  and hence  $y_{t+1} = \pm B_{t+1}$ , so

$$\begin{aligned}
 V_t(\mathbf{s}_t, \sigma_t^2) &= \min_{\hat{y}_{t+1}} \hat{y}_{t+1}^2 + 2B_{t+1} |\mathbf{x}_{t+1}^\top \mathbf{P}_{t+1} \mathbf{s}_t - \hat{y}_{t+1}| \\
 &\quad + \mathbf{x}_{t+1}^\top \mathbf{P}_{t+1} \mathbf{x}_{t+1} B^2 + \mathbf{s}_t^\top \mathbf{P}_{t+1} \mathbf{s}_t - \sigma_t^2 + \gamma_{t+1}.
 \end{aligned}$$

- ▶ Backwards induction gives

$$\begin{aligned}
 V_t(\mathbf{s}_t, \sigma_t^2) &:= \min_{\hat{y}_{t+1}} \max_{y_{t+1}} (\hat{y}_{t+1} - y_{t+1})^2 \\
 &\quad + V_{t+1}(\mathbf{s}_t + y_{t+1} \mathbf{x}_{t+1}, \sigma_t^2 + y_{t+1}^2), \\
 &= \min_{\hat{y}_{t+1}} \max_{y_{t+1}} (\hat{y}_{t+1} - y_{t+1})^2 \\
 &\quad + (\mathbf{s}_t + y_{t+1} \mathbf{x}_{t+1})^\top \mathbf{P}_{t+1} (\mathbf{s}_t + y_{t+1} \mathbf{x}_{t+1}) \\
 &\quad - (\sigma_t^2 + y_{t+1}^2) + \gamma_{t+1}
 \end{aligned}$$

- ▶ This is convex in  $y_{t+1}$  and hence  $y_{t+1} = \pm B_{t+1}$ , so

$$\begin{aligned}
 V_t(\mathbf{s}_t, \sigma_t^2) &= \min_{\hat{y}_{t+1}} \hat{y}_{t+1}^2 + 2B_{t+1} |\mathbf{x}_{t+1}^\top \mathbf{P}_{t+1} \mathbf{s}_t - \hat{y}_{t+1}| \\
 &\quad + \mathbf{x}_{t+1}^\top \mathbf{P}_{t+1} \mathbf{x}_{t+1} B_{t+1}^2 + \mathbf{s}_t^\top \mathbf{P}_{t+1} \mathbf{s}_t - \sigma_t^2 + \gamma_{t+1}.
 \end{aligned}$$

- ▶ If  $|\mathbf{x}_{t+1}^\top \mathbf{P}_{t+1} \mathbf{s}_t| \leq B_{t+1}$ , setting subgradient to 0 yields

$$\hat{y}_{t+1} = \mathbf{x}_{t+1}^\top \mathbf{P}_{t+1} \mathbf{s}_t$$

- ▶ Plugging in this  $\hat{y}_{t+1}$ , we get

$$V_t(\mathbf{s}_t, \sigma^2_t) = \mathbf{s}_t^\top (P_{t+1} \mathbf{x}_{t+1} \mathbf{x}_{t+1}^\top P_{t+1} + P_{t+1}) \mathbf{s}_t - \sigma^2_t + \gamma_{t+1} + B_{t+1}^2 \mathbf{x}_{t+1}^\top P_{t+1} \mathbf{x}_{t+1}$$

- ▶ Value is

$$V_0(\mathbf{0}, 0) = \gamma_0 = \sum_{t=1}^T B_t^2 \mathbf{x}_t^\top P_t \mathbf{x}_t$$

- ▶ Plugging in this  $\hat{y}_{t+1}$ , we get

$$V_t(\mathbf{s}_t, \sigma_t^2) = \mathbf{s}_t^\top \overbrace{\left( \mathbf{P}_{t+1} \mathbf{x}_{t+1} \mathbf{x}_{t+1}^\top \mathbf{P}_{t+1} + \mathbf{P}_{t+1} \right)}{:= \mathbf{P}_t} \mathbf{s}_t - \sigma_t^2 + \gamma_{t+1} + B_{t+1}^2 \mathbf{x}_{t+1}^\top \mathbf{P}_{t+1} \mathbf{x}_{t+1}$$

- ▶ Value is

$$V_0(\mathbf{0}, 0) = \gamma_0 = \sum_{t=1}^T B_t^2 \mathbf{x}_t^\top \mathbf{P}_t \mathbf{x}_t$$

- ▶ Plugging in this  $\hat{y}_{t+1}$ , we get

$$V_t(\mathbf{s}_t, \sigma^2_t) = \mathbf{s}_t^\top \left( \overbrace{P_{t+1} \mathbf{x}_{t+1} \mathbf{x}_{t+1}^\top P_{t+1} + P_{t+1}}^{:= P_t} \right) \mathbf{s}_t - \sigma^2_t + \underbrace{\gamma_{t+1} + B_{t+1}^2 \mathbf{x}_{t+1}^\top P_{t+1} \mathbf{x}_{t+1}}_{:= \gamma_t}$$

- ▶ Value is

$$V_0(\mathbf{0}, 0) = \gamma_0 = \sum_{t=1}^T B_t^2 \mathbf{x}_t^\top P_t \mathbf{x}_t$$

## Theorem

The strategy

$$\hat{y}_{t+1} = \mathbf{x}_{t+1}^\top \mathbf{P}_{t+1} \mathbf{s}_t, \quad (\text{MM})$$

is minimax optimal and the value-to-go is

$$V_t(\mathbf{s}_t, \sigma_t^2) = \mathbf{s}_t^\top \mathbf{P}_t \mathbf{s}_t - \sigma_t^2 + \gamma_t,$$

with coefficients

$$\mathbf{P}_T = \left( \sum_{t=1}^T \mathbf{x}_t \mathbf{x}_t^\top \right)^{-1}, \quad \mathbf{P}_t = \mathbf{P}_{t+1} + \mathbf{P}_{t+1} \mathbf{x}_{t+1} \mathbf{x}_{t+1}^\top \mathbf{P}_{t+1},$$
$$\gamma_T = 0, \quad \gamma_t = \gamma_{t+1} + B_{t+1}^2 \mathbf{x}_{t+1}^\top \mathbf{P}_{t+1} \mathbf{x}_{t+1},$$

provided the box constraints  $|\mathbf{x}_{t+1}^\top \mathbf{P}_{t+1} \mathbf{s}_t| \leq B_{t+1}$  hold.

## Alternate form of $P_t$

- ▶  $P_t$  has a nice interpretation as an augmented least squares prediction

$$P_t^{-1} = \underbrace{\sum_{q=1}^t \mathbf{x}_q \mathbf{x}_q^\top}_{\text{least squares}} + \underbrace{\sum_{q=t+1}^T \frac{\mathbf{x}_q^\top P_q \mathbf{x}_q}{1 + \mathbf{x}_q^\top P_q \mathbf{x}_q} \mathbf{x}_q \mathbf{x}_q^\top}_{\text{re-weighted future instances}}.$$



## Alternate form of $P_t$

- ▶  $P_t$  has a nice interpretation as an augmented least squares prediction

$$P_t^{-1} = \underbrace{\sum_{q=1}^t \mathbf{x}_q \mathbf{x}_q^T}_{\text{least squares}} + \underbrace{\sum_{q=t+1}^T \frac{\mathbf{x}_q^T P_q \mathbf{x}_q}{1 + \mathbf{x}_q^T P_q \mathbf{x}_q} \mathbf{x}_q \mathbf{x}_q^T}_{\text{re-weighted future instances}}.$$

- ▶ Accounts for future covariates
- ▶ Scale invariant
- ▶ Unlike ridge etc., data dependent regularization

# Regret

If the budgets and covariates are compatible, i.e. we have

$$B_t \geq \sum_{\tau=1}^{t-1} |\mathbf{x}_t^\top \mathbf{P}_t \mathbf{x}_\tau| B_\tau,$$

then the minimax regret is

$$\sum_{t=1}^T B_t^2 \mathbf{x}_t^\top \mathbf{P}_t \mathbf{x}_t$$

and the maximin probability distribution for  $\mathbf{y}_{t+1}$  puts weight  $1/2 \pm \mathbf{x}_{t+1}^\top \mathbf{P}_{t+1} \mathbf{s}_t / (2B_{t+1})$  on  $\pm B_{t+1}$ .

## Section 3

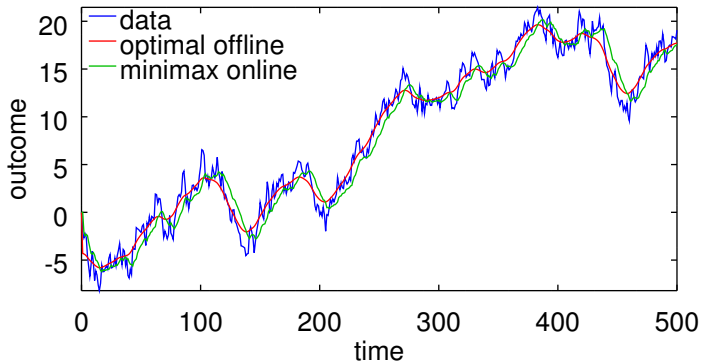
### Tracking

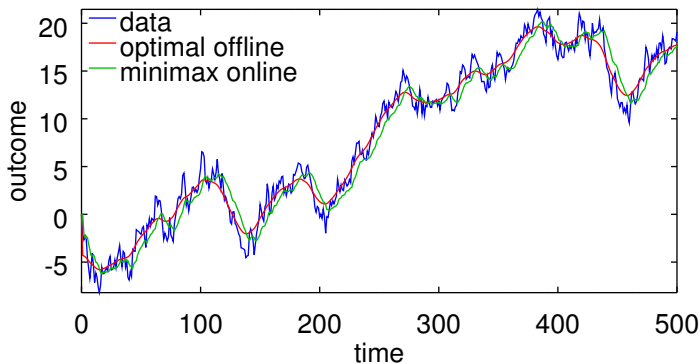
# Time series prediction protocol (with Koolen, Bartlett, Abbasi-Yadkori)

Fix a convex set  $\mathcal{C}$ , length  $T$ , regularization parameter  $\lambda_T$ .  
For each round  $t = 1, \dots, T$ ,

- ▶ We play  $\mathbf{a}_t \in \mathcal{C}$
- ▶ Nature reveals  $\mathbf{y}_t \in \mathcal{C}$
- ▶ We incur loss  $\ell(\mathbf{a}_t, \mathbf{y}_t) := \|\mathbf{a}_t - \mathbf{y}_t\|^2$
- ▶ Regret:

$$\underbrace{\sum_{t=1}^T \|\mathbf{a}_t - \mathbf{y}_t\|^2}_{\text{Our loss}} - \min_{\hat{\mathbf{a}}_1, \dots, \hat{\mathbf{a}}_T} \left\{ \underbrace{\sum_{t=1}^T \|\hat{\mathbf{a}}_t - \mathbf{y}_t\|^2}_{\text{Loss of Comparator}} + \underbrace{\lambda_T \sum_{t=1}^{T+1} \|\hat{\mathbf{a}}_t - \hat{\mathbf{a}}_{t-1}\|^2}_{\text{Comparator Complexity}} \right\}$$





Let

$Y_t = [y_1 \cdots y_t]$  and  $\hat{A} = [\hat{a}_1 \cdots \hat{a}_T]$ . For  $v_t \in \mathbb{R}^t$  and  $K \succeq 0$ ,

Data domain  $\|Y_t v_t\| \leq 1$  e.g.  $\|y_t\| \leq 1$

Complexity  $\text{tr}(K \hat{A}^\top \hat{A})$  e.g.  $\sum_{t=1}^{T+1} \|\hat{a}_t - \hat{a}_{t-1}\|^2$

## Backwards induction

Histories are  $\mathbf{Y}_t = [y_1 \cdots y_t]$ .

Offline Problem:  $\hat{\mathbf{A}} = \mathbf{Y}_T(\mathbf{I} + \lambda_T \mathbf{K})^{-1}$  and value

$$V_T(\mathbf{Y}_T) = -L^* = -\text{tr}(\mathbf{Y}_T(\mathbf{I} - (\mathbf{I} + \lambda_T \mathbf{K})^{-1})\mathbf{Y}_T^\top)$$

with recursion

$$V_{t-1}(\mathbf{Y}_{t-1}) = \min_{\mathbf{a}_t} \max_{y_t: \|\mathbf{Y}_t \mathbf{v}_t\| \leq 1} \|\mathbf{a}_t - y_t\|^2 + V_t(\mathbf{Y}_t).$$

So far, just a bit more complicated than before.

# Behavior of backwards induction solution

## Theorem

If  $\|\mathbf{b}\| \leq 1$ , then the minimax problem

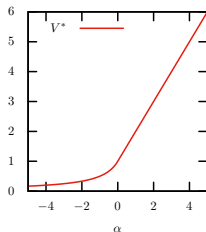
$$V^* := \min_{\mathbf{a}} \max_{\mathbf{y}: \|\mathbf{y}\| \leq 1} \|\mathbf{a} - \mathbf{y}\|^2 + (\alpha - 1)\|\mathbf{y}\|^2 + 2\mathbf{b}^T \mathbf{y}$$

has value and minimizer

$$V^* = \begin{cases} \frac{\|\mathbf{b}\|^2}{1-\alpha} & \text{if } \alpha \leq 0, \\ \|\mathbf{b}\|^2 + \alpha & \text{if } \alpha \geq 0, \end{cases} \quad \text{and} \quad \mathbf{a} = \begin{cases} \frac{\mathbf{b}}{1-\alpha} & \text{if } \alpha \leq 0, \\ \mathbf{b} & \text{if } \alpha \geq 0. \end{cases}$$

Non-trivial induction:

- ▶ Curvature of optimization can switch between rounds
- ▶ Yet can pre-compute beforehand





## Minimax solution

Input:  $T, K, \lambda_T, v_1, \dots, v_T$

Using:

- ▶ single-shot game solution, and
- ▶ lots of matrix identities

Output: matrices  $R_t = \begin{pmatrix} A_t & b_t \\ b_t^\top & c_t \end{pmatrix}$

strategy  $a_t = X_{t-1} \begin{cases} \frac{b_t}{1-c_t} & \text{if } c_t \leq 0, \\ b_t - c_t v_t^{<t} & \text{if } c_t \geq 0. \end{cases}$

## Theorem

Under a (typical) no clipping condition on  $Y_T$ ,

$$V(Y_t) = \text{tr}(Y_t(\mathbf{R}_t - \mathbf{I})Y_t^T) + \sum_{s=t+1}^T \max\{c_s, 0\}$$

and, in the vanilla case (norm bounded data, increments penalized),

$$V_t = \Theta\left(\frac{T}{\sqrt{1 + \lambda_T}}\right).$$

Section 4

Conclusion

- ▶ Minimax algorithms can be computationally efficient with enough structure, e.g.
  - ▶ Normalized Maximum likelihood that is Bayesian
  - ▶ Certain square losses
- ▶ Exploited the fact that saddle point problems with square loss are nice
- ▶ Can we characterize the class of functions that are closed w.r.t. the backwards induction operator?

## Section 5

Extra slides

## Ball game maximin

The maximin strategy plays two unit length vectors with

$$\Pr \left( y = \mathbf{a}_{\perp} \pm \sqrt{1 - \mathbf{a}_{\perp}^T \mathbf{a}_{\perp}} \mathbf{v}_{\max} \right) = \frac{1}{2} \pm \frac{\mathbf{a}_{\parallel}^T \mathbf{v}_{\max}}{2\sqrt{1 - \mathbf{a}_{\perp}^T \mathbf{a}_{\perp}}},$$

where  $\lambda_{\max}$  and  $\mathbf{v}_{\max}$  correspond to the largest eigenvalue of  $\mathbf{A}_{t+1}$  and  $\mathbf{a}_{\perp}$  and  $\mathbf{a}_{\parallel}$  are the components of  $\mathbf{a}^*$  perpendicular and parallel to  $\mathbf{v}_{\max}$ .

## Tracking: second order $K$

- ▶ Computation: if  $K$  and  $v_t$  are banded then  $R_t^{-1}$  is sparse
- ▶ Here we *imposed* data bound  $\|Y_t v_t\| \leq 1$ . In the paper we show that the minimax strategy guarantees an *adaptive* bound scaling with  $\|Y_t v_t\|$ .
- ▶ A second order smoothness version of  $K$  gives complicated  $c_t$

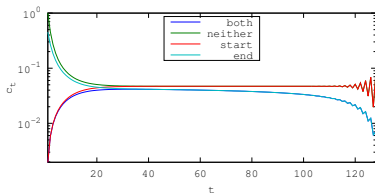


Figure:  $v_t = e_t - e_{t-1}$

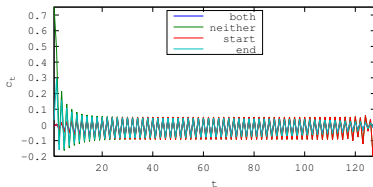


Figure:  $v_t = e_t - 2e_{t-1} + e_{t-2}$

## Ellipse

Fix a budget  $R \geq 0$ , and consider label sequences

$$\mathcal{Y}_R := \left\{ y_1, \dots, y_T \in \mathbb{R} : \sum_{t=1}^T y_t^2 \mathbf{x}_t^\top \mathbf{P}_t \mathbf{x}_t = R \right\}$$

We show that (MM) is minimax for this set.

In fact, the regret of (MM) *equals*

$$\mathcal{R}_T = \sum_{t=1}^T y_t^2 \mathbf{x}_t^\top \mathbf{P}_t \mathbf{x}_t.$$

This means that this algorithm has two very special properties. First, it is a *strong equalizer* in the sense that it suffers the same regret on all  $2^T$  sign-flips of the labels. And second, it is *adaptive* to the complexity  $R$  of the labels.