

High-Confidence Predictions Under Adversarial Uncertainty

Andrew Drucker

IAS

Setting: prediction on binary sequences

$$x = (x_1, x_2, x_3, \dots) \in \{0, 1\}^\omega$$

- Bits of x revealed sequentially.
- **Goal:** make some nontrivial prediction about unseen bits of sequence x , given bits seen so far.



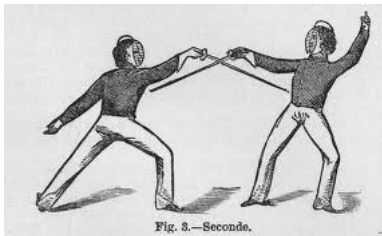
Setting: prediction on binary sequences

$$x = (x_1, x_2, x_3, \dots)$$

- **Question:** What kinds of assumptions on x are needed to make interesting predictions?
- **Our message:** Surprisingly weak ones.

Modeling questions

- Prediction: a game between the Predictor and Nature.
- What kind of opponent is Nature?



Probabilistic models

$$x = (x_1, x_2, x_3, \dots)$$

$$x \sim \mathcal{D},$$

where \mathcal{D} is some known probability distribution.

- **Problem:** how to choose correct \mathcal{D} for realistic applications?

Classes of assumptions

$$x = (x_1, x_2, x_3, \dots)$$

- **Adversarial models:**

$$x \in A,$$

where $A \subseteq \{0, 1\}^\omega$ is some known set.

- Interested in worst-case performance.
- These assumptions can be quite “safe” ...
- Our focus today.

Prior work on adversarial prediction

Gales and fractal dimension

[Lutz '03; Athreya, Hitchcock, Lutz, Mayordomo '07]

- Gales: a class of betting strategies, to bet on unseen bits of $x \in A$.
- Goal: reach a fortune of ∞ , on any $x \in A$.
- The “handicap” we need can be related to measures of fractal dimension for A ...

Prior work on adversarial prediction

Ignorant forecasting

- What is the chance of rain tomorrow?



- Basic test of a meteorologist: “calibration.”
- If governing distribution \mathcal{D} is known, easy to achieve with Bayes' rule...
- But: calibration can also be achieved by an ignorant forecaster! [**Foster, Vohra '98**]

Prior work on adversarial prediction

$$x = (x_1, x_2, x_3, \dots)$$

- These works' goal: long-term, overall predictive success.
- Our focus: make a single prediction with high confidence.

0-prediction

- **Our main scenario:** want to predict a single 0 among the bits of x .

(We lose if prediction fails OR if we wait forever!)

- **Interpretation:** choose a time to “safely” perform some action;

$[x_t = 0]$ means “time t is safe.”



Possible assumptions

- **ϵ -biased arrivals assumption:** bits of x independent, with

$$\Pr[x_t = 1] = \epsilon.$$

- Best strategy succeeds with prob. $1 - \epsilon$.

Possible assumptions

- Very strong assumption...
- **Idea** (not new): study adversarial “relaxations” of ϵ -biased model.

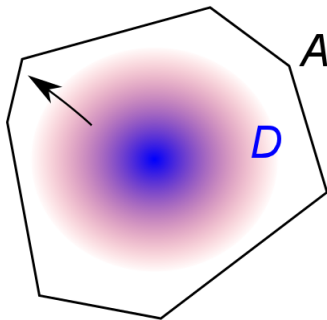
Possible assumptions

- Very strong assumption...
- **Idea** (not new): study adversarial “relaxations” of ϵ -biased model.



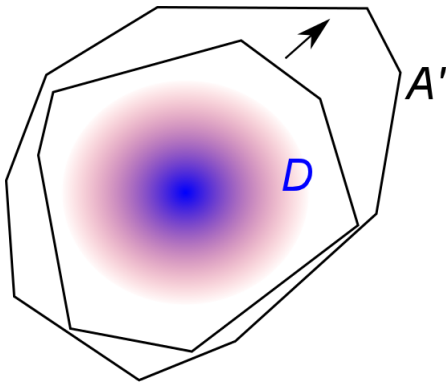
Possible assumptions

- Very strong assumption...
- **Idea** (not new): study adversarial “relaxations” of ϵ -biased model.



Possible assumptions

- Very strong assumption...
- **Idea** (not new): study adversarial “relaxations” of ϵ -biased model.



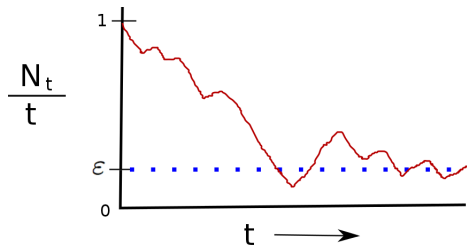
Possible assumptions

- Let

$$N_t := x_1 + \dots + x_t.$$

- ε -sparsity assumption:** say that x is ε -sparse if

$$\limsup_{t \rightarrow \infty} N_t/t \leq \varepsilon.$$



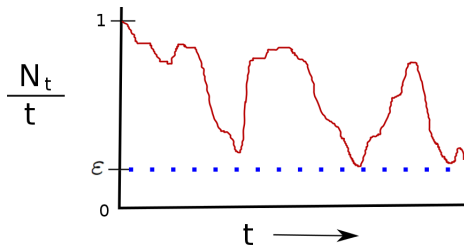
Possible assumptions

- Let

$$N_t := x_1 + \dots + x_t.$$

- ε -weak sparsity assumption:** say that x is ε -weakly sparse if

$$\liminf_{t \rightarrow \infty} N_t/t \leq \varepsilon.$$



Our main result

Theorem

For any $\varepsilon, \gamma > 0$, there is a (randomized) 0-prediction strategy $\mathcal{S}_{\varepsilon, \gamma}$ that succeeds with prob. $\geq 1 - \varepsilon - \gamma$, on any ε -weakly sparse sequence.

- Can do nearly as well as under ε -biased arrivals!

Our main result

Theorem

For any $\varepsilon, \gamma > 0$, there is a (randomized) 0-prediction strategy $\mathcal{S}_{\varepsilon, \gamma}$ that succeeds with prob. $\geq 1 - \varepsilon - \gamma$, on any ε -weakly sparse sequence.

- Can do nearly as well as under ε -biased arrivals!
- (Adversary's sequence gets fixed before randomness in $\mathcal{S}_{\varepsilon, \gamma} \dots$)

Proof ideas

- Divide sequence into “epochs:”

$$x = \underbrace{0}_{E_1} \underbrace{1110}_{E_2} \underbrace{010110000}_{E_3} \dots$$

- (r -th epoch of length $K_r = \Theta(r^2)$.)
- Run a separate 0-prediction algorithm for each individual epoch.

Proof ideas

$$x = \underbrace{0}_{E_1} \underbrace{1110}_{E_2} \underbrace{010110000}_{E_3} \dots$$

- **Easy claim:** x is ε -weakly sparse



∃ a subsequence of “**nice**” epochs,
whose 1-densities are at most $\varepsilon + \gamma/3$.

Let $\varepsilon' = \varepsilon + \gamma/2$.

Proof ideas

$$x = \underbrace{0}_{E_1} \underbrace{1110}_{E_2} \underbrace{010110000}_{E_3} \dots$$

Idea: give an algorithm \mathcal{S} with the properties:

- 1 Makes a 0-prediction with noticeable prob. on each nice epoch;
- 2 On every epoch,

$$\Pr \left[\text{true prediction} \right] \geq \left(\frac{1-\epsilon'}{\epsilon'} \right) \cdot \Pr \left[\text{false prediction} \right].$$

(Would achieve our goal!)

Proof ideas

$$\Pr \left[\text{true prediction} \right] \stackrel{?}{\geq} \left(\frac{1-\epsilon'}{\epsilon'} \right) \cdot \Pr \left[\text{false prediction} \right]$$

- Whoops—can't achieve this!
- **Modified goal:** an upper bound

$$\left(\frac{1-\epsilon'}{\epsilon'} \right) \cdot \Pr \left[\text{false prediction} \right] - \Pr \left[\text{true prediction} \right] \leq \text{(small)}$$

Proof ideas

$$\Pr \left[\text{true prediction} \right] \stackrel{?}{\geq} \left(\frac{1-\epsilon'}{\epsilon'} \right) \cdot \Pr \left[\text{false prediction} \right]$$

- Whoops—can't achieve this!
- **Modified goal:** an upper bound

$$\left(\frac{1-\epsilon'}{\epsilon'} \right) \cdot \Pr \left[\text{false prediction} \right] - \Pr \left[\text{true prediction} \right] \leq O(1/|K_r|)$$

Proof ideas

- During the r -th epoch, alg. maintains a stack of “chips” (initially empty);



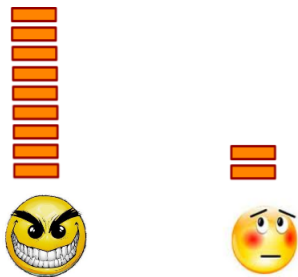
stack's height



algorithm's “**courage**” to predict next bit of x will be a 0.

Proof ideas

- During the r -th epoch, alg. maintains a stack of “chips” (initially empty);



stack's height



algorithm's “**courage**” to predict next bit of x will be a 0.

Stack dynamics

- Assume

$$\varepsilon' = \frac{p}{d} = 1 - \frac{q}{d}.$$

- Observe a 0: add p “courage chips.”
- Observe a 1: remove q chips.

e.g., $p = 1$:

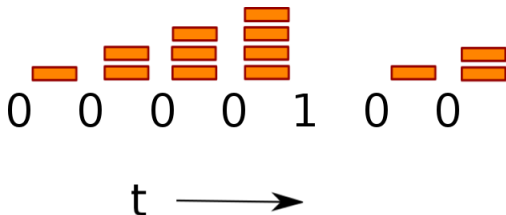
Stack dynamics

- Assume

$$\varepsilon' = \frac{p}{d} = 1 - \frac{q}{d}.$$

- Observe a 0: add p "courage chips."
- Observe a 1: remove q chips.

e.g., $p = 1$:



Making predictions

- Let $H_t =$ stack height after observing first t bits of r -th epoch.
- Overall algorithm for epoch r :**
 - Choose t^* uniformly from $\{1, 2, \dots, K_r\}$;
 - Observe first $t^* - 1$ bits;

0 0 1 0 0 1 **?** ? ? ? ?
t*

- Predict a 0 on step t^* with probability

$$\frac{H_{t^*-1}}{d \cdot K_r},$$

else make no prediction this epoch.

Analysis ideas

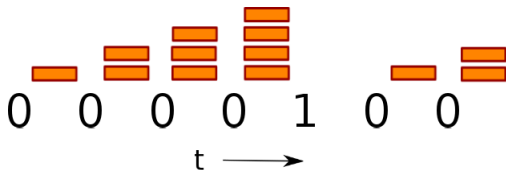
- 1 If fraction of 1s in r -th epoch is $\leq \varepsilon + \gamma/3 < p/d$, a 0-prediction is made in epoch r with $\Omega(1)$ prob.



\Rightarrow Eventually (in some epoch), a prediction is made.

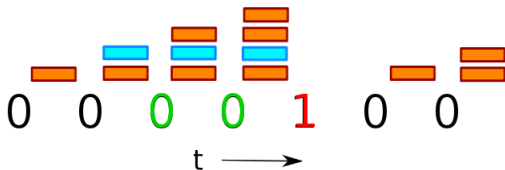
Analysis ideas

- 2 To compare odds of correct and incorrect 0-predictions, analyze each chip's contribution.



Analysis ideas

- 2 To compare odds of correct and incorrect 0-predictions, analyze each chip's contribution.



Analysis ideas

- **Intuition:**

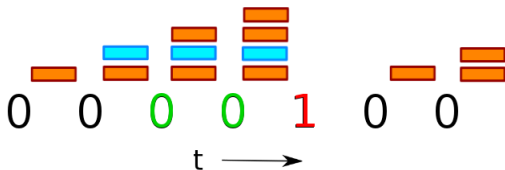
- If a chip remains on the stack long enough, fraction of 1s while it's on is $\lesssim p/d = \epsilon'$.

GOOD! (Contributes mostly to successful predictions.)

- Total contribution to failure probability of other (“bad”) chips is small.
- We can analyze all chips in a simple, unified way...

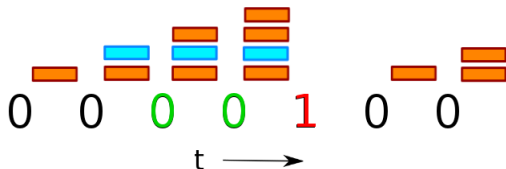
Analysis ideas

- Fix attention to a chip c on input x .



- Let $\text{zeros}(c)$ ($\text{ones}(c)$) denote the number of zeros (ones) appearing after steps where c is on the stack.

Analysis ideas

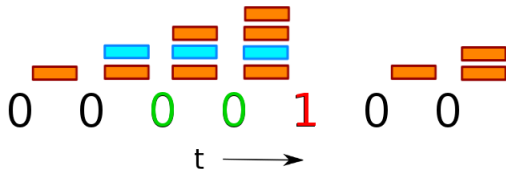


- Chip c 's contribution to **success** and **failure** probabilities:

$$\begin{array}{ccc} \downarrow & & \downarrow \\ \frac{\text{zeros}(c)}{d \cdot K_r^2}, & & \frac{\text{ones}(c)}{d \cdot K_r^2}. \end{array}$$

- To compare: show that $\text{zeros}(c)$ and $\text{ones}(c)$ obey a linear inequality...

Analysis ideas



Claim:

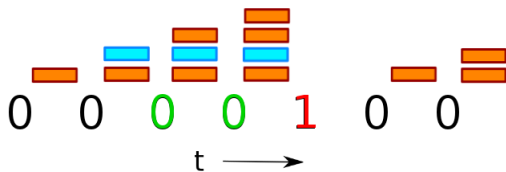
$$q \cdot \text{ones}(c) - p \cdot \text{zeros}(c) \leq q.$$

Proof: LHS bounded by the net loss in stack height between first appearance of c and (possible) removal...

c is removed along with $\leq q$ other chips!

Let's sum over all c ...

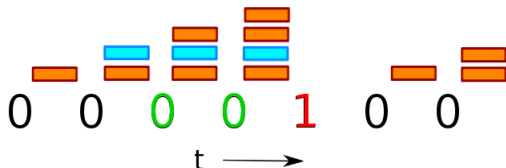
Analysis ideas



Summing over all c (at most $p \cdot K_r$ chips total):

$$q \cdot \sum_c \text{ones}(c) - p \cdot \sum_c \text{zeros}(c) \leq pqK_r.$$

Analysis ideas



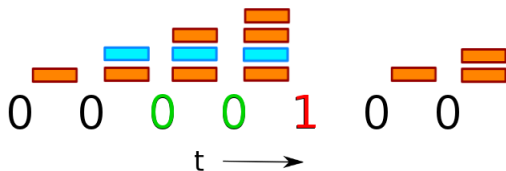
Summing over all c (at most $p \cdot K_r$ chips total):

$$(q/p) \cdot \frac{\sum_c \text{ones}(c)}{d \cdot K_r^2} - \sum_c \frac{\text{zeros}(c)}{d \cdot K_r^2} \leq \frac{q}{dK_r}.$$

$$(q/p) \cdot \Pr[\text{failure in epoch } r] - \Pr[\text{success in epoch } r]$$

$$\leq O(K_r^{-1}) = O(r^{-2}).$$

Analysis ideas



Summing over all c (at most $p \cdot K_r$ chips total):

$$(q/p) \cdot \frac{\sum_c \text{ones}(c)}{d \cdot K_r^2} - \sum_c \frac{\text{zeros}(c)}{d \cdot K_r^2} \leq \frac{q}{dK_r}.$$

$$\left(\frac{1 - \epsilon'}{\epsilon'} \right) \cdot \Pr[\text{failure in epoch } r] - \Pr[\text{success in epoch } r] \\ \leq O(K_r^{-1}) = O(r^{-2}).$$

Q.E.D.

Also in the paper

- Bit prediction for broader classes of assumptions:
- E.g., predict a bit (0 or 1), under the assumption that a certain word appears only rarely.
- General statement involves finite automata.

Also in the paper

- Also: high-confidence predictions under no assumptions on x !

Ignorant interval-forecasting

- Sequence $x \in \{0, 1\}^\omega$: now completely arbitrary.
- **Goal:** predict the fraction of 1s in an unseen interval, with high accuracy and high confidence.

(Huh?)

Ignorant interval-forecasting

- **Our “hook”**—we get to choose the position and size of the prediction-interval.
- **Interval-forecaster alg.:** makes a prediction of form:

“A p fraction of the next N bits will be 1s.”

Ignorant interval-forecasting

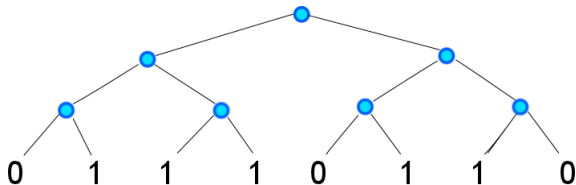
Theorem

For any $\epsilon, \delta > 0$, there is a ignorant interval-forecaster $\mathcal{S}_{\epsilon, \delta}$ that is accurate to $\pm\epsilon$, with success probability $1 - \delta$.

Runtime of $\mathcal{S}_{\epsilon, \delta}$ is finite: $= 2^{O(\epsilon^{-2}\delta^{-1})}$.

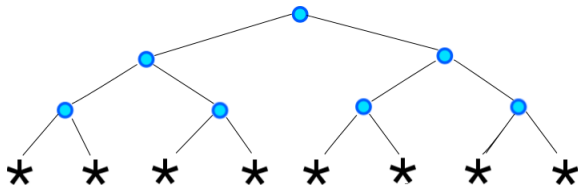
The approach

- Consider $x \in \{0, 1\}^{2^n}$, $n = \lceil 4/(\varepsilon^2 \delta) \rceil$.
- Arrange bits of x on leaves of a binary tree \mathcal{T} .



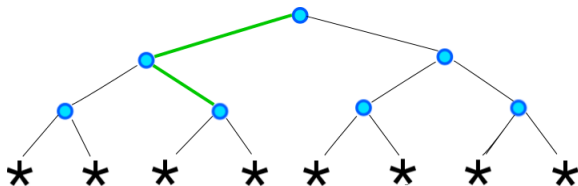
The approach

Forecasting algorithm:



The approach

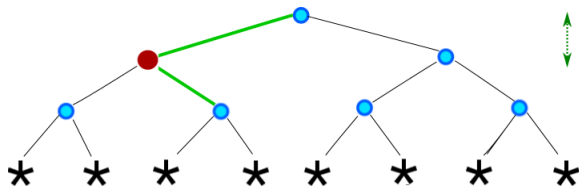
Forecasting algorithm:



- 1 Choose a random walk \mathcal{W} from root in \mathcal{T} of length $n - 1$;

The approach

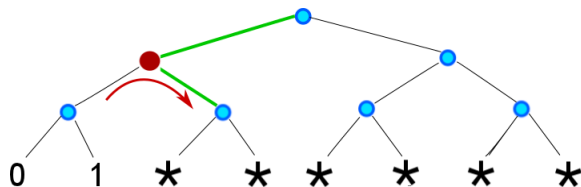
Forecasting algorithm:



- 1 Choose a random walk \mathcal{W} from root in \mathcal{T} of length $n - 1$;
- 2 Pick a uniform $t^* \in \{0, 1, \dots, n - 1\}$, and select t^{*th} vertex along \mathcal{W} ;

The approach

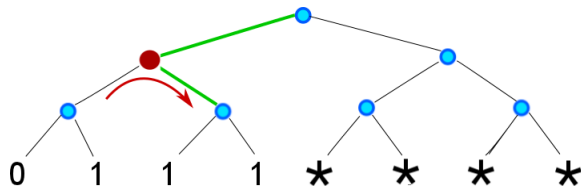
Forecasting algorithm:



- 1 Choose a random walk \mathcal{W} from root in \mathcal{T} of length $n - 1$;
- 2 Pick a uniform $t^* \in \{0, 1, \dots, n - 1\}$, and select t^{*th} vertex along \mathcal{W} ;
- 3 **Predict** that
(fraction of 1s in right subtree) =
(fraction of 1s in left subtree).

The approach

Forecasting algorithm:



- 1 Choose a random walk \mathcal{W} from root in \mathcal{T} of length $n - 1$;
- 2 Pick a uniform $t^* \in \{0, 1, \dots, n - 1\}$, and select t^{*th} vertex along \mathcal{W} ;
- 3 **Predict** that
(fraction of 1s in right subtree) =
(fraction of 1s in left subtree).

Analysis idea

- For $0 \leq t \leq n$ let

$$X_t \in [0, 1]$$

denote the fraction of 1s below the t -th step vertex.

- **Fact:** For any fixed bit-sequence x ,

$$X_0, X_1, \dots, X_n$$

is a **martingale**, from which:

$$\mathbb{E}[(X_{t+1} - X_t)(X_{s+1} - X_s)] = 0,$$

for all $s < t$. Thus:

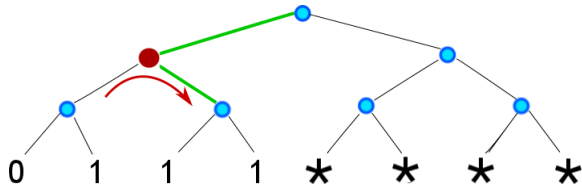
$$1 \geq \mathbb{E}[(X_n - X_0)^2] = \sum_{0 \leq t < n} \mathbb{E}[(X_{t+1} - X_t)^2].$$

Analysis idea

$$\sum_{0 \leq t < n} \mathbb{E}[(X_{t+1} - X_t)^2] \leq 1$$

- $(X_{t+1} - X_t)^2$ small $\implies t$ is a good choice for t^* !

(i.e., left and right subtrees have similar 1-densities).



- So w.h.p. over walk \mathcal{W} , most choices for t^* are good!
Q.E.D.

Questions

- Characterize the sets $A \subset \{0, 1\}^\omega$ for which confident 0-prediction is possible?

Connection with fractal dimension, à la (Lutz et al.)?

Questions

- For which distributions \mathcal{D} on $\{0, 1\}^\infty$ can we extend to a “supporting set” A , preserving easiness of prediction?

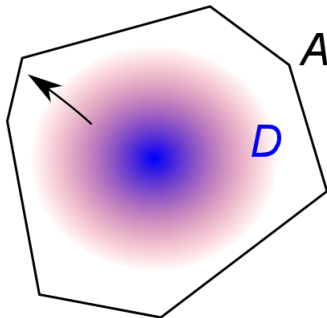
Questions

- For which distributions \mathcal{D} on $\{0, 1\}^\infty$ can we extend to a “supporting set” A , preserving easiness of prediction?



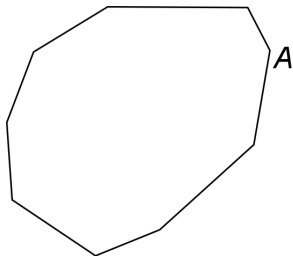
Questions

- For which distributions \mathcal{D} on $\{0, 1\}^\infty$ can we extend to a “supporting set” A , preserving easiness of prediction?



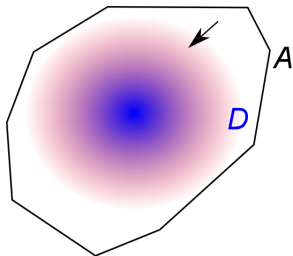
Questions

- Is there a minimax theorem for 0-prediction?
- Hard set A for 0-prediction \Rightarrow hard distribution D over A ?



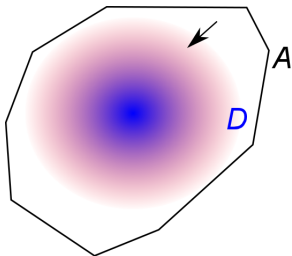
Questions

- Is there a minimax theorem for 0-prediction?
- Hard set A for 0-prediction \Rightarrow hard distribution D over A ?



Questions

- Is there a minimax theorem for 0-prediction?
- Hard set A for 0-prediction \Rightarrow hard distribution D over A ?



- Would give alternate (non-constructive) proof of our main result...
- More examples of surprisingly confident prediction?

Thanks!