

On-Line Learning

Nicolò Cesa-Bianchi

Università degli Studi di Milano

November 3, 2005



- Theory of repeated games
(Hannan, 1956; Blackwell, 1956)
- Compression of individual sequences
(Lempel and Ziv, 1976)
- Gambling and portfolio selection
(Cover, 1965 and 1991)
- Pattern classification
(Novikov, 1962; Littlestone, 1989)

Unifying framework

Prediction with expert advice



- 1 Prediction with expert advice
- 2 Connections with game theory
- 3 Learning with linear experts
- 4 The Perceptron algorithm and its extensions
- 5 Mistake bounds
- 6 Online learning with kernels
- 7 From mistake bounds to risk bounds



Binary prediction

- A **forecaster** predicts a binary sequence one bit at the time
- At each step $t = 1, 2, \dots$ the forecaster predicts the t -th bit knowing the previous $t - 1$ bits

0100010110?...

- After the prediction is made, the t -th bit is observed and the forecaster finds out whether a mistake was made

Goal

Bound the number of prediction mistakes without making any statistical assumptions on the way the data sequence is generated



The role of experts

- Want a nonstatistical framework where **good** forecasters can be distinguished from **bad** forecasters
- Any forecaster must use some map of the form

past observations \rightarrow predictions

- For each forecaster, there exists a bit sequence on which a mistake is made at each step

Competitive analysis

Compare the performance of the forecaster to that of a set of *reference forecasters* (**experts**)



A simple example

Forecaster competes against three experts on sequence 1101

	$t = 1$	$t = 2$	$t = 3$	$t = 4$	Mistakes
Expert 1	1	1	1	1	$M_1 = 1$
Expert 2	0	1	1	0	$M_2 = 3$
Expert 3	1	0	1	0	$M_3 = 3$
Forecaster	1	0	1	1	$M = 2$
Bit sequence	1	1	0	1	

Goal (refined)

Predict each sequence almost as well as the best expert for that sequence



A more general prediction model

- Predict an unknown sequence $y_1, y_2, \dots \in \mathcal{Y}$ (outcome space)
- Predictions \hat{p} are chosen from \mathcal{X} (decision space)
- Forecasters are scored with their cumulative loss

$$\ell(\hat{p}_1, y_1) + \ell(\hat{p}_2, y_2) + \dots$$

where $\ell : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ is a **loss function**

Example

- **Zero-one loss:** $\mathcal{X} = \mathcal{Y} = \{0, 1\}$ and $\ell(\hat{p}, y) = \mathbb{I}_{\{\hat{p} \neq y\}}$
- **Quadratic loss:** $\mathcal{X} = \mathcal{Y} = [0, 1]$ and $\ell(\hat{p}, y) = (\hat{p} - y)^2$
- **Absolute loss:** $\mathcal{X} = [0, 1], \mathcal{Y} = \{0, 1\}$ and $\ell(\hat{p}, y) = |\hat{p} - y|$



On-line prediction with expert advice

Measure performance relatively to a set of N experts

At each step $t = 1, 2, \dots$

- 1 Get predictions (advice) $f_{1,t}, \dots, f_{N,t} \in \mathcal{X}$ of the experts
- 2 Compute prediction $\hat{p}_t \in \mathcal{X}$
- 3 Outcome $y_t \in \mathcal{Y}$ is revealed
- 4 Forecaster incurs loss $\ell(\hat{p}_t, y_t)$ and each expert i incurs loss $\ell(f_{i,t}, y_t)$

Note

Experts are viewed as **abstract entities**, generating predictions in an unspecified way



$$r_{i,t} = \ell(\hat{p}_t, y_t) - \ell(f_{i,t}, y_t)$$

$$\hat{L}_n = \sum_{t=1}^n \ell(\hat{p}_t, y_t) \quad L_{i,n} = \sum_{t=1}^n \ell(f_{i,t}, y_t)$$

$$R_{i,n} = \sum_{t=1}^n r_{i,t} = \hat{L}_n - L_{i,n}$$

$$\max_{i=1,\dots,N} R_{i,n} = \hat{L}_n - \min_{i=1,\dots,N} L_{i,n}$$

We want to design **consistent** forecasters, i.e. such that

$$\lim_{n \rightarrow \infty} \frac{1}{n} \left(\max_{i=1,\dots,N} R_{i,n} \right) = 0$$

for any sequence of outcomes and all choices of expert advice



Weighted average forecasters

- Assume decision space \mathcal{X} is a **convex subset** of a linear space

$$x, x' \in \mathcal{X} \implies \alpha x + (1 - \alpha)x' \in \mathcal{X} \quad \text{for all } 0 \leq \alpha \leq 1$$

- If $R_{i,t-1}$ is big, then we should predict more like expert i

$$\hat{p}_t = \frac{\sum_{i=1}^N \mu(R_{i,t-1}) f_{i,t}}{\sum_{j=1}^N \mu(R_{j,t-1})}$$

where μ is some positive monotone increasing function

- This is the **weighted average** forecaster



Convex loss functions

- Assume loss function $\ell(\mathbf{x}, \mathbf{y})$ is convex in its first argument
- Then, the prediction $\hat{\mathbf{p}}_t$ satisfies

$$\ell(\hat{\mathbf{p}}_t, \mathbf{y}_t) = \ell\left(\frac{\sum_{i=1}^N \mu(R_{i,t-1}) \mathbf{f}_{i,t}}{\sum_{j=1}^N \mu(R_{j,t-1})}, \mathbf{y}_t\right) \leq \frac{\sum_{i=1}^N \mu(R_{i,t-1}) \ell(\mathbf{f}_{i,t}, \mathbf{y}_t)}{\sum_{j=1}^N \mu(R_{j,t-1})}$$

- Using $r_{i,t} = \ell(\hat{\mathbf{p}}_t, \mathbf{y}_t) - \ell(\mathbf{f}_{i,t}, \mathbf{y}_t)$ and rearranging we obtain that, irrespective to \mathbf{y}_t ,

$$\sum_{i=1}^N r_{i,t} \mu(R_{i,t-1}) \leq 0$$



Potential-based forecasters

- Choose $\mu = \phi'$
where $\phi : \mathbb{R} \rightarrow \mathbb{R}$ is s.t. $\phi, \phi' \geq 0$ and ϕ'' exists
- Weighted average forecaster is then

$$\hat{p}_t = \frac{\sum_{i=1}^N \phi'(R_{i,t-1}) f_{i,t}}{\sum_{j=1}^N \phi'(R_{j,t-1})}$$

Definition

Potential function $\Phi : \mathbb{R}^N \rightarrow \mathbb{R}$

$$\Phi(\mathbf{R}) = \psi \left(\sum_{i=1}^N \phi(R_i) \right)$$

where $\psi : \mathbb{R} \rightarrow \mathbb{R}$ is such that $\psi \geq 0, \psi' > 0, \psi'' \leq 0$



Blackwell condition

Then the prediction at time t becomes

$$\hat{p}_t = \frac{\sum_{i=1}^N \nabla \Phi(R_{i,t-1})_i f_{i,t}}{\sum_{j=1}^N \nabla \Phi(R_{j,t-1})_j}$$

And the condition

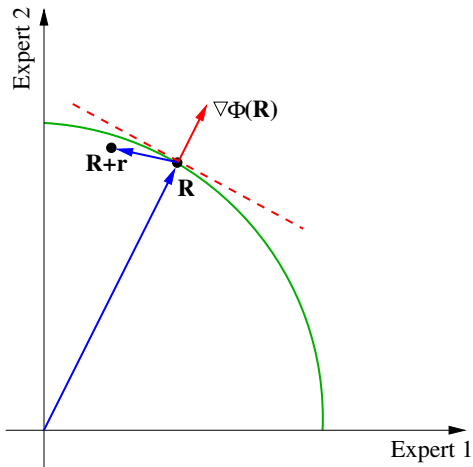
$$\sum_{i=1}^N r_{i,t} \mu(R_{i,t-1}) \leq 0$$

gets rewritten as

$$\nabla \Phi(\mathbf{R}_{t-1})^\top \mathbf{r}_t \leq 0 \quad (\text{Blackwell condition})$$



Gradient descent interpretation



Proving regret bounds

Since ϕ is monotone increasing

$$\begin{aligned}\psi\left(\phi\left(\max_{i=1,\dots,N} R_{i,n}\right)\right) &= \psi\left(\max_{i=1,\dots,N} \phi(R_{i,n})\right) \\ &\leq \psi\left(\sum_{i=1}^N \phi(R_{i,n})\right) = \Phi(\mathbf{R}_n)\end{aligned}$$

Assuming ψ is also invertible,

$$\max_{i=1,\dots,N} R_{i,n} \leq \phi^{-1}\psi^{-1}\left(\Phi(\mathbf{R}_n)\right)$$

So, a bound on $\Phi(\mathbf{R}_n)$ implies a bound on $\max_{i=1,\dots,N} R_{i,n}$



Proof of the potential bound

We bound the increment $\Phi(\mathbf{R}_t) - \Phi(\mathbf{R}_{t-1})$ of the potential by taking a **linear approximation** of $\Phi(\mathbf{R}_t)$ around $\Phi(\mathbf{R}_{t-1})$

$$\begin{aligned}\Phi(\mathbf{R}_t) &= \Phi(\mathbf{R}_{t-1} + \mathbf{r}_t) \\ &= \Phi(\mathbf{R}_{t-1}) + \nabla\Phi(\mathbf{R}_{t-1})^\top \mathbf{r}_t \\ &\quad + \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \left. \frac{\partial^2 \Phi(\mathbf{u})}{\partial u_i \partial u_j} \right|_{\mathbf{u}=\xi} r_{i,t} r_{j,t} \quad (\text{for some } \xi \in \mathbb{R}^N) \\ &\leq \Phi(\mathbf{R}_{t-1}) + \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \left. \frac{\partial^2 \Phi(\mathbf{u})}{\partial u_i \partial u_j} \right|_{\mathbf{u}=\xi} r_{i,t} r_{j,t}\end{aligned}$$



As $\Phi(\mathbf{u}) = \psi \left(\sum_{k=1}^N \phi(u_k) \right) = \psi(\Sigma(\mathbf{u}))$, we have that

$$\frac{\partial \Phi(\mathbf{u})}{\partial u_i} = \psi'(\Sigma(\mathbf{u})) \phi'(u_i)$$

For $i \neq j$,

$$\frac{\partial^2 \Phi(\mathbf{u})}{\partial u_i \partial u_j} = \psi''(\Sigma(\mathbf{u})) \phi'(u_i) \phi'(u_j)$$

For $i = j$,

$$\frac{\partial^2 \Phi(\mathbf{u})}{\partial u_i^2} = \psi''(\Sigma(\mathbf{u})) \phi'(u_i)^2 + \psi'(\Sigma(\mathbf{u})) \phi''(u_i)$$



Hence,

$$\begin{aligned} & \sum_{i=1}^N \sum_{j=1}^N \frac{\partial^2 \Phi(\mathbf{u})}{\partial u_i \partial u_j} \Big|_{\mathbf{u}=\xi} r_{i,t} r_{j,t} \\ &= \psi''(\Sigma(\xi)) \sum_{i=1}^N \sum_{j=1}^N \phi'(\xi_i) \phi'(\xi_j) r_{i,t} r_{j,t} \\ & \quad + \psi'(\Sigma(\xi)) \sum_{i=1}^N \phi''(\xi_i) r_{i,t}^2 \\ &= \psi''(\Sigma(\xi)) \left(\sum_{i=1}^N \phi'(\xi_i) r_{i,t} \right)^2 + \psi'(\Sigma(\xi)) \sum_{i=1}^N \phi''(\xi_i) r_{i,t}^2 \\ &\leq \psi'(\Sigma(\xi)) \sum_{i=1}^N \phi''(\xi_i) r_{i,t}^2 \end{aligned}$$



- We have proven that

$$\Phi(\mathbf{R}_t) \leq \Phi(\mathbf{R}_{t-1}) + \underbrace{\frac{1}{2} \max_{\mathbf{u}} \psi' \left(\sum_{k=1}^N \phi(\mathbf{u}_k) \right) \sum_{i=1}^N \phi''(\mathbf{u}_i) r_{i,t}^2}_{C(\mathbf{r}_t)}$$

- $C(\mathbf{r}_t)$ bounds **Taylor error**

- By iterating, we get $\Phi(\mathbf{R}_n) \leq \Phi(\mathbf{0}) + \frac{1}{2} \sum_{t=1}^n C(\mathbf{r}_t)$

- This holds for any forecaster satisfying $\nabla \Phi(\mathbf{R}_{t-1})^\top \mathbf{r}_t \leq 0$
(e.g., weighted average forecaster for convex losses)



Polynomial potential

Assume: Loss ℓ is convex and takes values in $[0, 1]$

- Potential function

$$\Phi_p(\mathbf{R}) = \left(\sum_{i=1}^N (R_i)_+^p \right)^{2/p} = \|\mathbf{R}\|_p^2 \quad \text{for } p \geq 2$$

- Prediction

$$\hat{p}_t = \frac{\sum_{i=1}^N \phi'(R_{i,t-1}) f_{i,t}}{\sum_{j=1}^N \phi'(R_{j,t-1})} = \frac{\sum_{i=1}^N (R_{i,t-1})_+^{p-1} f_{i,t}}{\sum_{j=1}^N (R_{j,t-1})_+^{p-1}}$$

- Taylor error bound: $C(\mathbf{r}_t) = (p-1) \|\mathbf{r}_t\|_p^2 \leq (p-1)N^{2/p}$

- Bound on regret: $\max_{i=1, \dots, N} R_{i,n} \leq \sqrt{n(p-1)N^{2/p}}$



Comments on polynomial potential

- The polynomial forecaster is consistent with rate $O(1/\sqrt{n})$
- There is a trade-off for the choice of p in the bound

$$\max_{i=1,\dots,N} R_{i,N} \leq \sqrt{n(p-1)N^{2/p}}$$

- Choosing $p = 2 \ln N$ yields

$$\max_{i=1,\dots,N} R_{i,N} \leq \sqrt{(2e)n \ln N}$$

- Is this the **best possible bound** in term of n and N ?



Exponential potential

Assume: Loss ℓ is convex and takes values in $[0, 1]$

- Potential function

$$\Phi_\eta(\mathbf{R}) = \frac{1}{\eta} \ln \left(\sum_{i=1}^N e^{\eta R_i} \right) \quad \text{for } \eta > 0$$

- Prediction:

$$\hat{p}_t = \frac{\sum_{i=1}^N e^{\eta(\hat{L}_{t-1} - L_{i,t-1})} f_{i,t}}{\sum_{j=1}^N e^{\eta(\hat{L}_{t-1} - L_{j,t-1})}} = \frac{\sum_{i=1}^N e^{-\eta L_{i,t-1}} f_{i,t}}{\sum_{j=1}^N e^{-\eta L_{j,t-1}}}$$

- Taylor error bound: $C(\mathbf{r}_t) = \eta \max_{i=1, \dots, N} r_{i,t}^2 \leq \eta$

- Bound on regret: $\max_{i=1, \dots, N} R_{i,n} \leq \frac{\ln N}{\eta} + \frac{\eta}{2} n$



Tuned exponential forecaster

- Choosing $\eta = \sqrt{2(\ln N)/n}$ gets

$$\max_{i=1,\dots,N} R_{i,n} \leq \sqrt{2n \ln N}$$

- Better constant than polynomial potential but not consistent (**horizon-dependent** tuning)
- Disregarding consistency, is the constant $\sqrt{2}$ optimal?



Lower bound

- Best lower bound is for the **absolute loss** $\ell(p, y) = |p - y|$ ($p \in [0, 1]$ and $y \in \{0, 1\}$)

$$\max_{y_1, \dots, y_n} \frac{1}{n} \left(\max_{i=1, \dots, N} R_{i,n} \right) = (1 - o(1)) \sqrt{\frac{n}{2} \ln N}$$

for **any** forecasting strategy $o(1) \rightarrow 0$ for $N, n \rightarrow \infty$

- Refining the analysis of the exponential potential we get the **optimal constant**

$$\max_{i=1, \dots, N} R_{i,N} \leq \sqrt{\frac{n}{2} \ln N}$$

- But tuning of η is still horizon-dependent



More sophisticated forecasters

- Using a **time-dependent** tuning $\eta_t = \sqrt{8(\ln N)/t}$ we get

$$\max_{i=1,\dots,N} R_{i,N} \leq \sqrt{2n \ln N} + \sqrt{\frac{\ln N}{8}}$$

- Consistent, but suboptimal leading constant

- Loss-dependent** tuning $\eta_t = \sqrt{c \frac{\ln N}{L_{t-1}^*}}$ gives

$$\max_{i=1,\dots,N} R_{i,N} \leq 2 \sqrt{2L_n^* \ln N} + O(\ln N)$$

$$L_t^* = \min_{i=1,\dots,N} L_{i,t}$$



- Bound for forecasters satisfying **Blackwell condition**

$$\Phi(\mathbf{R}_n) \leq \Phi(\mathbf{0}) + \frac{1}{2} \sum_{t=1}^n C(\mathbf{r}_t)$$

- Polynomial potential with $p = 2 \ln N$

$$\max_{i=1,\dots,N} R_{i,n} \leq \sqrt{(2e)n \ln N}$$

- Exponential potential with time-varying parameter

$$\max_{i=1,\dots,N} R_{i,n} \leq \sqrt{2n \ln N} + \sqrt{\frac{\ln N}{8}}$$



The greedy forecaster for the exponential potential

- Predicts by minimizing the **increase of regret**

$$\begin{aligned}\hat{p}_t &= \operatorname{argmin}_{\hat{p} \in \mathcal{X}} \sup_{y_t \in \mathcal{Y}} \Phi(\mathbf{R}_{t-1} + \mathbf{r}_t) \\ &= \operatorname{argmin}_{\hat{p} \in \mathcal{X}} \sup_{y_t \in \mathcal{Y}} \left(\ell(\hat{p}, y_t) + \frac{1}{\eta} \ln \sum_{i=1}^N e^{-\eta L_{i,t}} \right)\end{aligned}$$

- Not necessarily a weighted average forecaster

Definition

A loss is **mixable** if for some η^* and for all $t = 1, 2, \dots$ the greedy forecaster satisfies

$$\ell(\hat{p}_t, y_t) \leq -\frac{1}{\eta} \ln \frac{\sum_{i=1}^N e^{-\eta^* L_{i,t}}}{\sum_{j=1}^N e^{-\eta^* L_{j,t-1}}}$$



Properties of mixable losses

- If a loss is mixable for η^* , then $\Phi_{\eta^*}(\mathbf{R}_n) \leq \Phi_{\eta^*}(\mathbf{0})$
- In fact, we have

$$\hat{L}_n \leq -\frac{1}{\eta^*} \ln \frac{\sum_{j=1}^N e^{-\eta^* L_{j,n}}}{N} \leq L_{i,n} + \frac{\ln N}{\eta^*} \quad \text{for any } i$$

- This immediately implies

$$\max_{i=1,\dots,N} R_{i,n} \leq \frac{1}{\eta^*} \ln N$$

- A **constant** regret bound!



Which losses are mixable?

- Assume $0 \leq \hat{p}, y \leq 1$
- **Square loss** $\ell(\hat{p}, y) = (\hat{p} - y)^2$ is mixable for $\eta^* = 2$
(greedy \neq weighted average)
- **Relative entropy loss** $\ell(\hat{p}, y) = y \ln \frac{y}{\hat{p}} + (1 - y) \ln \frac{1 - y}{1 - \hat{p}}$ is mixable for $\eta^* = 1$ (greedy = weighted average)
- **Absolute loss** $\ell(\hat{p}, y) = |\hat{p} - y|$ is not mixable for any $\eta^* > 0$



Zero-sum games

$N \times M$ **known loss matrix** with entries $0 \leq \ell(i, y) \leq 1$

$\ell(1,1)$	$\ell(1,2)$...
$\ell(2,1)$	$\ell(2,2)$...
\vdots	\vdots	\ddots

Row player has N actions

Column player has
 M actions

- Players independently draw actions I (with law P) and Y (with law Q)
- Row player suffers loss $\ell(I, Y)$ (= gain of column player)
- **Value of the game**

$$\min_P \max_Q \sum_{i,y} \ell(i,y)P(i)Q(y) = \max_Q \min_P \sum_{i,y} \ell(i,y)P(i)Q(y)$$



Repeated zero-sum games

- Game is played repeatedly
- At round $t = 1, 2, \dots$ players draw actions (I_t, Y_t) which may depend on **past realizations** of (I_s, Y_s) , $s = 1, \dots, t - 1$
- **Regret after n plays**

$$R_{i,n} = \sum_{t=1}^n \ell(I_t, Y_t) - \sum_{t=1}^n \ell(i, Y_t)$$

- Regret is small when empirical distribution of row player's actions performs not much worse than any fixed action



Definition

A row player is **Hannan consistent** if

$$\limsup_{n \rightarrow \infty} \max_{i=1, \dots, N} \frac{R_{i,n}}{n} = 0 \quad \text{with probability 1}$$

irrespective to what column player does

- H.C. implies $\limsup_{n \rightarrow \infty} \frac{1}{n} \sum_{t=1}^n \ell(I_t, Y_t)$ is at most the game value
- **Note:** Row player can beat the game value if column player is suboptimal



Randomized forecasting

At each round t :

- 1 **forecaster** announces distribution P_t over $\{1, \dots, N\}$
 - 2 **adversary** picks $y_t \in \mathcal{Y}$ and **forecaster** draws $I_t \sim P_t$
 - 3 I_t and y_t are both revealed
- We may apply the exp. weighted average forecaster

$$P_t(i) \propto \exp\left(-\eta_t \sum_{s=1}^{t-1} \ell(i, y_s)\right) \quad \eta_t = \sqrt{\frac{8 \ln N}{t}}$$

- This achieves Hannan consistency

$$\max_{i=1, \dots, N} \frac{R_{i,n}}{n} \leq \left(\sqrt{2 \frac{\ln N}{n}} + o(1) \right) \quad \text{w.h.p.}$$



Label efficient prediction

- What if forecaster queries only m out of n adversary's actions?
- For any fixed n , we get

$$\max_{i=1,\dots,N} \frac{R_{i,n}}{n} \leq c \sqrt{\frac{\ln N}{m}} \quad \text{w.h.p.}$$

- Optimal to within constants
- A query rate slightly faster than

$$\frac{(\ln n)(\ln \ln n)}{n}$$

is sufficient for Hannan consistency



Partial monitoring

- What if forecaster observes **signal** $h(I_t, y_t)$ instead of y_t ?
- **Feedback matrix** H

Example

- Forecaster's action $I_t \in \{1, 2, \dots, N\}$ is the price at which a product sold online is offered to t -th customer
- Adversary's action $y_t \in \{1, 2, \dots, N\}$ is maximum price at which t -th customer is willing to buy the product
- Feedback matrix is

$$h(I_t, y_t) = \begin{cases} \text{SOLD} & \text{if } I_t \leq y_t \\ \text{NOT SOLD} & \text{otherwise} \end{cases}$$



Partial monitoring and Hannan consistency

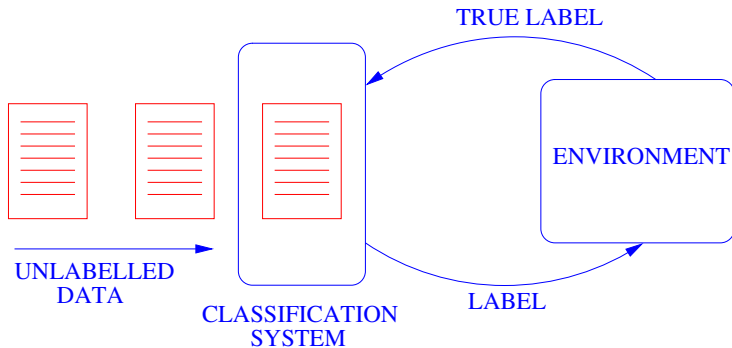
- Assume \mathcal{Y} is finite
- L is the loss matrix, H is any (legal) encoding of the feedback matrix
- Sufficient (and almost necessary) condition:

$$L = KH \quad \text{for some matrix } K$$

- Then per-round regret vanishes at rate $n^{-1/3}$
- Optimal, but worse than $n^{-1/2}$ of the **full monitoring** scenario
- Can still get $n^{-1/2}$ in special cases (e.g., nonstochastic bandits)



Binary pattern classification



Pattern classification model

- Data instances encoded as vectors $\mathbf{x}_t \in [0, 1]^d$
- A binary label $y_t \in \{-1, 1\}$ expresses some property of \mathbf{x}_t

Example

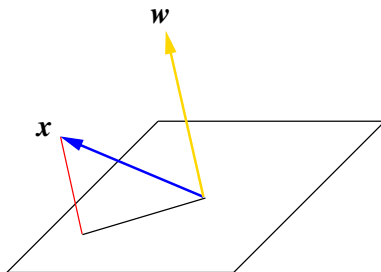
Given today's closing prices \mathbf{x}_t , predict whether tomorrow market index will increase ($y_t = 1$)



Linear classifiers

Linear classification

predict with $\hat{p}_t = \text{SGN}(\mathbf{w}_{t-1}^\top \mathbf{x}_t)$ $\mathbf{w}_{t-1} \in \mathbb{R}^d$



Linear classifiers (cont.)

If $\hat{p}_t \neq y_t$ then **mistake at step t**

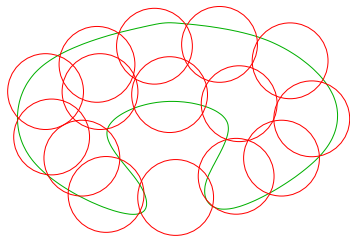
Goal

On any arbitrary sequence $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots$ perform not much worse than the **best fixed linear classifier**

Experts = all linear classifiers



Direct application of experts' framework



- Consider the class \mathcal{F} of all linear classifiers $\hat{p}_t = \text{SGN}(\mathbf{u}^\top \mathbf{x}_t)$ for $\mathbf{u} \in \mathbb{R}^d$ with $\|\mathbf{u}\|$ bounded
- A covering of \mathcal{F} has size **exponential** in d
- Running the weighted average forecaster on the covering requires managing an exponential number of weights



A reduction to prediction with expert advice

- Allocate d experts F_1, \dots, F_d
- On instance $\mathbf{x}_t = (x_{t,1}, \dots, x_{t,d})$ expert F_j predicts $x_{t,j}$

Regret

$$\mathbf{r}_t = y_t \mathbf{x}_t \mathbb{I}_{\{\hat{p}_t \neq y_t\}}$$

We write \mathbf{m}_t to denote $\mathbb{I}_{\{\hat{p}_t \neq y_t\}}$



A reduction (cont.)

- Unnormalized weighted average forecaster for binary classification

$$\mathbf{w}_{t-1} = \nabla \Phi(\mathbf{R}_{t-1}) \quad \hat{p}_t = \text{SGN}(\mathbf{w}_{t-1}^\top \mathbf{x}_t)$$

- To apply previous results we need **Blackwell condition** $\mathbf{w}_{t-1}^\top \mathbf{r}_t \leq 0$ to hold
- Indeed,

$$\mathbf{w}_{t-1}^\top \mathbf{r}_t = y_t \mathbf{w}_{t-1}^\top \mathbf{x}_t \mathbb{I}_{\{\hat{p}_t \neq y_t\}} = \begin{cases} 0 & \text{if } \mathbb{I}_{\{\hat{p}_t \neq y_t\}} = 0 \\ < 0 & \text{otherwise} \end{cases}$$

since $\mathbb{I}_{\{\hat{p}_t \neq y_t\}} = 1$ iff $\text{SGN}(\mathbf{w}_{t-1}^\top \mathbf{x}_t) \neq y_t$



Polynomial potential

- Recall

$$\Phi_p(\mathbf{R}) = \left(\sum_{i=1}^N R_i^p \right)^{2/p} = \|\mathbf{R}\|_p^2 \quad \text{for } p \geq 2$$

Before we had $\Phi_p(\mathbf{R}) = \|(\mathbf{R})_+\|_p^2$

Now $(\cdot)_+$ dropped for technical reasons

- Applying the main result

$$\|\mathbf{R}_n\|_p^2 \leq \frac{p-1}{2} \sum_{t=1}^n \|\mathbf{r}_t\|_p^2 \leq \frac{p-1}{2} \left(\max_t \|\mathbf{x}_t\|_p \right)^2 \sum_{t=1}^n m_t$$

- Next we lower bound $\|\mathbf{R}_n\|_p = \sqrt{\Phi_p(\mathbf{R}_n)}$



Polynomial potential (cont.)

- For any linear classifier of parameter \mathbf{u}

$$\begin{aligned}\|\mathbf{R}_n\|_p &\geq \mathbf{R}_n^\top \frac{\mathbf{u}}{\|\mathbf{u}\|_q} \quad (\text{by Hölder's inequality}) \\ &= \mathbf{R}_{n-1}^\top \frac{\mathbf{u}}{\|\mathbf{u}\|_q} + m_n \frac{y_n \mathbf{u}^\top \mathbf{x}_n}{\|\mathbf{u}\|_q} \\ &\geq \mathbf{R}_{n-1}^\top \frac{\mathbf{u}}{\|\mathbf{u}\|_q} + m_n \frac{1 - d_n(\mathbf{u})}{\|\mathbf{u}\|_q}\end{aligned}$$

where $d_n(\mathbf{u}) = (1 - y_n \mathbf{u}^\top \mathbf{x}_n)_+$ is the **hinge loss** of \mathbf{u}

- Iterating we get

$$\|\mathbf{R}_n\|_p \geq \sum_{t=1}^n \frac{m_t - d_t(\mathbf{u})}{\|\mathbf{u}\|_q}$$



A general mistake bound

- Piecing together upper and lower bounds:

$$\sum_{t=1}^n \frac{m_t - d_t(\mathbf{u})}{\|\mathbf{u}\|_q} \leq \|\mathbf{R}_n\|_p \leq \sqrt{\frac{(p-1)X_p^2}{2} \sum_{t=1}^n m_t}$$

- Solving for $M_n = \sum_{t=1}^n m_t$ and overapproximating

$$M_n - D_n(\mathbf{u}) \leq \frac{p-1}{2} (X_p \|\mathbf{u}\|_q)^2 + (X_p \|\mathbf{u}\|_q) \sqrt{\frac{p-1}{2} D_n(\mathbf{u})}$$

where $D_n(\mathbf{u}) = \sum_{t=1}^n d_t(\mathbf{u})$

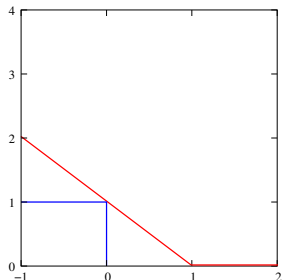
This holds for **all** $\mathbf{u} \in \mathbb{R}^d$ and **all** sequences $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots \in \mathbb{R}^d \times \{-1, +1\}$



The hinge loss

The hinge loss upper bounds the mistake indicator function

$$d_t = \left(1 - y_t \mathbf{u}^\top \mathbf{x}_t\right)_+ \geq \text{SGN}(\mathbf{u}^\top \mathbf{x}_t)$$



This “regret” is a bit unfair

$$M_n - \inf_{\mathbf{u}} D_n(\mathbf{u}) \leq M_n - \inf_{\mathbf{u}} M_n(\mathbf{u})$$



Formulation as an incremental algorithm

We want to express $\mathbf{w}_t = \nabla\Phi(\mathbf{R}_t)$ recursively as $\mathbf{w}_t = F(\mathbf{w}_{t-1})$

Definition

A potential $\Phi : \mathbb{R}^d \rightarrow \mathbb{R}$ is **Legendre** if Φ is strictly convex, differentiable, and has a convex domain (plus some additional technical requirements)

If a potential is Legendre, then $\nabla\Phi$ is **invertible**

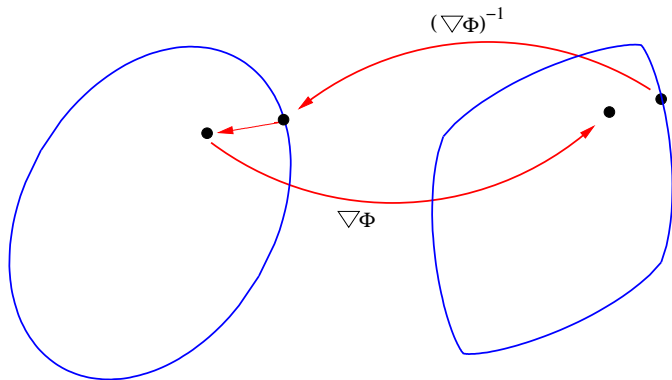
$$\mathbf{w}_t = \nabla\Phi(\mathbf{R}_t) = \nabla\Phi(\mathbf{R}_{t-1} + \mathbf{r}_t) = \nabla\Phi\left((\nabla\Phi)^{-1}(\mathbf{w}_{t-1}) + \mathbf{r}_t\right)$$

Incremental formulation

$$\mathbf{w}_t = \nabla\Phi\left((\nabla\Phi)^{-1}(\mathbf{w}_{t-1}) + \mathbf{y}_t \mathbf{x}_t \mathbb{I}_{\{\hat{\mathbf{p}}_t \neq \mathbf{y}_t\}}\right) \quad \text{update rule}$$



Incremental formulation (cont.)



$$\mathbf{w}_t = \nabla\Phi \left((\nabla\Phi)^{-1}(\mathbf{w}_{t-1}) + \mathbf{y}_t \mathbf{x}_t \mathbb{I}_{\{\hat{p}_t \neq y_t\}} \right)$$



Application to polynomial potential

- Polynomial potential $\Phi_p(\cdot) = \|\cdot\|_p^2$ is Legendre

$$\left(\nabla_{\frac{1}{2}}\|\mathbf{u}\|_p^2\right)_i = \frac{\text{SGN}(u_i) |u_i|^{p-1}}{\|\mathbf{u}\|_p^{p-2}} \quad \left(\nabla_{\frac{1}{2}}\|\mathbf{u}\|_p^2\right)^{-1} = \nabla_{\frac{1}{2}}\|\mathbf{u}\|_q^2$$

where q is such that $1/p + 1/q = 1$

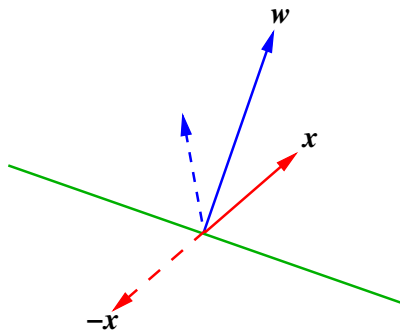
- When $p = 2$ we have $\nabla\Phi_2(\mathbf{R}) = \mathbf{R}$
- The incremental formulation then is simply

$$\mathbf{w}_t = \mathbf{w}_{t-1} + y_t \mathbf{x}_t \mathbb{I}_{\{\hat{p}_t \neq y_t\}}$$

the **Perceptron algorithm** (Rosenblatt, 1952)



The Perceptron algorithm



$$\mathbf{w}_t = \mathbf{w}_{t-1} + y_t \mathbf{x}_t \mathbb{I}_{\{\hat{p}_t \neq y_t\}}$$



Mistake bounds

We can specialize our previous bound by setting $p = 2$

$$M_n - D_n(\mathbf{u}) \leq (X_2 \|\mathbf{u}\|_2)^2 + (X_2 \|\mathbf{u}\|_2) \sqrt{D_n(\mathbf{u})}$$

Definition

A sequence $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$ is **linearly separable** with margin $\gamma > 0$ if there exists $\mathbf{u} \in \mathbb{R}^d$ with $\|\mathbf{u}\|_2 = 1$ such that

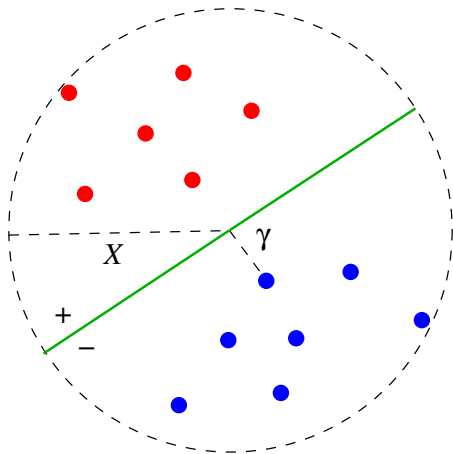
$$\min_{t=1, \dots, n} y_t \mathbf{u}^\top \mathbf{x}_t \geq \gamma$$

In this case we recover the **Perceptron convergence theorem**

$$M_n \leq \left(\frac{X_2}{\gamma} \right)^2$$



Linear separability



Exponential potential

- Recall

$$\Phi_{\eta}(\mathbf{R}) = \frac{1}{\eta} \ln \left(\sum_{i=1}^d e^{\eta R_i} \right)$$

- The weights have form

$$\mathbf{w}_t = \nabla \Phi_{\eta}(\mathbf{R}_t) = \frac{e^{\eta R_{i,n}}}{\sum_{k=1}^d e^{\eta R_{k,n}}}$$

- Note that \mathbf{w}_t belongs to the simplex in \mathbb{R}^d



Incremental formulation

- **Problem:** Φ_η is not strictly convex because $\nabla\Phi_\eta(\mathbf{R})$ is constant along the line $R_1 = R_2 = \dots = R_d$
- Thus, Φ is not Legendre and $(\nabla\Phi)^{-1}$ is not defined
- However, the potential $\Phi_{\text{exp}}(\mathbf{R}) = e^{R_1} + \dots + e^{R_d}$ is Legendre with

$$\begin{aligned}\nabla\Phi_{\text{exp}}(\mathbf{R}) &= (e^{R_1}, \dots, e^{R_d}) \\ (\nabla\Phi_{\text{exp}})^{-1}(\mathbf{w}') &= (\ln w'_1, \dots, \ln w'_d)\end{aligned}$$

- We get the following update rule

$$\begin{aligned}w'_{i,t} &= \left[\nabla\Phi_{\text{exp}} \left((\nabla\Phi_{\text{exp}})^{-1}(\mathbf{w}'_{t-1}) + \eta \mathbf{r}_{t-1} \right) \right]_i \\ &= \exp(\ln w'_{i,t-1} + \eta r_{i,t-1}) \\ &= w'_{i,t-1} e^{\eta r_{i,t-1}}\end{aligned}$$



Incremental formulation (cont.)

- Summarizing, we have

$$\begin{aligned}\mathbf{w}'_t &= \nabla \Phi_{\text{exp}} \left((\nabla \Phi_{\text{exp}})^{-1} (\mathbf{w}'_{t-1}) + \eta \mathbf{r}_{t-1} \right) \\ w_{i,t} &= \left[\nabla \Phi_{\eta}(\mathbf{R}_t) \right]_i = \frac{e^{\eta R_{i,t}}}{\sum_{k=1}^d e^{\eta R_{k,t}}} = \frac{w'_{i,t}}{\sum_{k=1}^d w'_{k,t}}\end{aligned}$$

- This is the **Winnow algorithm** (Littlestone, 1988)



Exponential potential (cont.)

- Applying the main result

$$\begin{aligned}\ln \Phi_{\eta}(\mathbf{R}_n) &\leq \frac{\ln d}{\eta} + \frac{\eta}{2} \sum_{t=1}^n \left(\max_{i=1, \dots, d} r_{i,t}^2 \right) \\ &\leq \frac{\ln d}{\eta} + \frac{\eta}{2} (X_{\infty})^2 M_n\end{aligned}$$

- As before, we lower bound $\ln \Phi_{\eta}(\mathbf{R}_n)$



- Recalling that \mathbf{w} belong to the simplex, we apply the **log-sum inequality**

$$\ln \sum_{i=1}^d v_i \geq \sum_{i=1}^d u_i \ln v_i + H(\mathbf{u}) \quad u_i, v_i \geq 0 \quad \sum_{i=1}^d u_i = 1$$

for any linear classifier \mathbf{u} in the simplex

- Applying this to the exponential potential we get

$$\ln \Phi_{\eta}(\mathbf{R}_n) = \frac{1}{\eta} \ln \sum_{i=1}^d e^{\eta R_{i,n}} \geq \mathbf{R}_n^T \mathbf{u} + H(\mathbf{u})$$

Dropping $H(\mathbf{u})$ and proceeding as for the polynomial potential we get

$$\ln \Phi_{\eta}(\mathbf{R}_n) \geq \sum_{t=1}^n (m_t - d_t(\mathbf{u})) = M_n - D_n(\mathbf{u})$$



- Piecing together upper and lower bound

$$M_n - D_n(\mathbf{u}) \leq \ln \Phi_\eta(\mathbf{R}_n) \leq \frac{\ln d}{\eta} + \frac{\eta}{2}(X_\infty)^2 M_n$$

- Solving for M_n

$$M_n \leq \frac{1}{1 - \eta(X_\infty)^2/2} \left(D_n(\mathbf{u}) + \frac{\ln d}{\eta} \right)$$

for any sequence $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots \in \mathbb{R}^d \times \{-1, +1\}$ and for any \mathbf{u} in the simplex.

- Same **tuning problem** as in prediction with expert advice
- Proper choice of $\eta = \eta(D_n, X_\infty, d)$ gives

$$M_n - D_n(\mathbf{u}) \leq X_\infty \sqrt{2D_n(\mathbf{u}) \ln d} + (1 + o(1))(X_\infty)^2 \ln d$$

where $o(1) \rightarrow 0$ for $D(\mathbf{u}) \rightarrow \infty$



Beyond the simplex

- Recall that, for the exponential potential, the difference $M_n - D_n(\mathbf{u})$ is bounded for all \mathbf{u} in the **simplex**
- If we consider \mathbf{u} in the **scaled simplex**
 $\{\mathbf{u} \in \mathbb{R}^d : u_i \geq 0, u_1 + \dots + u_d = U\}$ we get

$$M_n - D_n(\mathbf{u}) \leq (X_\infty U) \sqrt{2D_n(\mathbf{u}) \ln d} + (1 + o(1))(X_\infty U)^2 \ln d$$

- To remove the constraint $u_i \geq 0$ we can **transform the instances**

$$\mathbf{x} = (x_1, \dots, x_d) \mapsto (x_1, -x_1, \dots, x_d, -x_d)$$

at the price of replacing d by $2d$



Comparison between poly. and exp. potential

The two bounds on $M_n - D_n(\mathbf{u})$

$$\frac{p-1}{2} (X_p \|\mathbf{u}\|_q)^2 + (X_p \|\mathbf{u}\|_q) \sqrt{\frac{p-1}{2} D_n(\mathbf{u})}$$

$$(1 + o(1)) \ln(2d) (X_\infty \|\mathbf{u}\|_1)^2 + (X_\infty \|\mathbf{u}\|_1) \sqrt{2D_n(\mathbf{u}) \ln(2d)}$$

- Bound for exp. pot. assumes tuning (previous knowledge of X_∞ and choice of $\|\mathbf{u}\|_1$)
- Both bounds depend on pairs of **dual norms**: $\|\mathbf{x}\|_p \|\mathbf{u}\|_q$ vs. $\|\mathbf{x}\|_\infty \|\mathbf{u}\|_1$
- For $p \approx 2 \ln d$ the bounds are essentially equal



Comparison for spherical potential

- Consider a sequence $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2) \dots$ such that $\mathbf{x}_t \in \{-1, 1\}^d$ and $y_t = \text{SGN}(\mathbf{x}_{1,t})$
- Then $\mathbf{u} = (1, 0, \dots, 0)$ is an optimal classifier (no loss)
- Moreover,

$$(\|\mathbf{u}\|_2 X_2)^2 = d \quad \text{and} \quad (\|\mathbf{u}\|_1 X_\infty)^2 = 1$$

- Then

$$M_n \leq d \quad (\text{polynomial potential, } p = 2)$$

$$M_n \leq 4 \ln(2d) \quad (\text{exponential potential})$$

an exponential advantage (verified by experiments)

- Opposite situation when instances \mathbf{x}_t are **sparse** and best expert \mathbf{u} is **dense**



On-line learning with kernels

- Feature map $\phi : \mathbb{R}^d \rightarrow \text{RKHS}$
- Kernel $K(\mathbf{x}, \mathbf{x}') = \langle \phi(\mathbf{x}), \phi(\mathbf{x}') \rangle$
- Assume a linear algorithm learns \mathbf{w} such that

$$\mathbf{w} = \sum_i \alpha_i \mathbf{x}_{t_i}$$

- Then we can learn $\mathbf{w} = \sum_i \alpha_i \phi(\mathbf{x}_{t_i})$ in the RKHS because

$$\begin{aligned} \text{SGN}(\langle \mathbf{w}, \phi(\mathbf{x}) \rangle) &= \text{SGN} \left(\sum_i y_{t_i} \langle \phi(\mathbf{x}_{t_i}), \phi(\mathbf{x}) \rangle \right) \\ &= \text{SGN} \left(\sum_i y_{t_i} K(\mathbf{x}_{t_i}, \mathbf{x}) \right) \end{aligned}$$



Checking applicability of kernels

$$\text{Let } \mathbf{R}_t = \sum_t y_t \mathbf{x}_t \mathbb{I}_{\{\hat{p}_t \neq y_t\}}$$

- **Winnow** $w_{i,t} = \frac{e^{\eta R_{i,t}}}{\sum_{k=1}^d e^{\eta R_{k,t}}}$

- **p-norm Perceptron** $w_{i,t} = \frac{\text{SGN}(R_{i,t}) |R_{i,t}|^{p-1}}{\|\mathbf{R}_t\|_p^{p-2}}$

- **Perceptron** $\mathbf{w}_t = \mathbf{R}_t$

Perceptron's potential is spherical \rightarrow rotational invariance



Kernel Perceptron

Initialize $\mathcal{L} = \emptyset$

For $t = 1, 2, \dots$

- 1 Read next example (\mathbf{x}_t, y_t)
- 2 Compute prediction $\hat{p} = \text{SGN} \left(\sum_{k \in \mathcal{L}} y_k K(\mathbf{x}_k, \mathbf{x}_t) \right)$
- 3 If $\hat{p}_t \neq y_t$ then $\mathcal{L} \leftarrow \mathcal{L} \cup \{t\}$

Now mistake bounds are extended to the whole RKHS

$$M_n \leq D_n(f) + \left(\max_t K(\mathbf{x}_t, \mathbf{x}_t) \right) \|f\|^2 + \dots$$

for any f in the RKHS



- Linear classifiers $H(\mathbf{x}) = \text{sgn}(\mathbf{w}^\top \mathbf{x})$
- Examples (\mathbf{x}_t, y_t) are i.i.d. according to a fixed and unknown probability distribution on $\mathbb{R}^d \times \{-1, +1\}$
- $\text{risk}(H) = \mathbb{P}(H(\mathbf{x}) \neq y)$
- Learning algorithm

$$(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n) \longrightarrow \boxed{A} \longrightarrow \hat{H} : \mathbb{R}^d \rightarrow \{-1, +1\}$$

\hat{H} is (random) hypothesis output by learner



Data-dependent VC theory

- \mathcal{H} is set of classifiers from which \hat{H} is selected
- For all $H \in \mathcal{H}$

$$\begin{aligned} & \text{risk}(H) - \text{risk}_{\text{emp}}(H) \\ & \leq c_1 \sqrt{\text{risk}_{\text{emp}}(H) \frac{V_{\mathcal{H}} \ln n}{n}} + c_2 \frac{V_{\mathcal{H}} \ln n}{n} \quad \text{w.h.p.} \end{aligned}$$

- VC theory of generalization studies properties of \mathcal{H}
- We use on-line learning to study a small subclass of \mathcal{H} generated by the interaction with the training data



The ensemble of hypotheses

- Run an incremental learner on the training set
- Everytime $H(\mathbf{x}_t) \neq y_t$, H is changed by the update rule
- This process generates an **ensemble of classifiers**

$$H_0, H_1, \dots, H_n$$

Goals

- 1 Bound the **average risk of the ensemble** in terms of the size of the ensemble
- 2 Find an element of the ensemble whose risk is close to the ensemble average



Step 1: bound the average risk

The difference

$$\text{risk}(H_{t-1}) - \mathbb{I}_{\{H_{t-1}(\mathbf{x}_t) \neq y_t\}}$$

is a **martingale difference sequence** because

$$\mathbb{E} \left[\text{risk}(H_{t-1}) - \mathbb{I}_{\{H_{t-1}(\mathbf{x}_t) \neq y_t\}} \mid (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_{t-1}, y_{t-1}) \right] = 0$$

The associated martingale is

$$\begin{aligned} & \sum_{t=1}^n \left(\text{risk}(H_{t-1}) - \mathbb{I}_{\{H_{t-1}(\mathbf{x}_t) \neq y_t\}} \right) \\ & \iff \underbrace{\frac{1}{n} \sum_{t=1}^n \text{risk}(H_{t-1})}_{\text{average risk}} - \underbrace{\frac{1}{n} \sum_{t=1}^n \mathbb{I}_{\{H_{t-1}(\mathbf{x}_t) \neq y_t\}}}_{\text{fraction of mistakes}} \end{aligned}$$



Bernstein's bound

If Z_1, Z_2, \dots is a **martingale difference sequence** with increments bounded by 1 and

$$V_n = \sum_{t=1}^n \mathbb{E} [Z_t^2 \mid Z_1, \dots, Z_{t-1}]$$

then for all $S, K > 0$

$$\mathbb{P} \left(\sum_{t=1}^n Z_t \geq S, \quad V_n \leq K \right) \leq \exp \left(-\frac{S^2}{2(S/3 + K)} \right)$$



Application of Bernstein's bound

Since $0 \leq \mathbb{I}_{\{H(\mathbf{x}) \neq y\}} \leq 1$,

$$\begin{aligned} \text{VAR} \left[\mathbb{I}_{\{H_{t-1}(\mathbf{x}_t), y_t\}} \mid (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_{t-1}, y_{t-1}) \right] \\ \leq \mathbb{E} \left[\text{risk}(H_{t-1}) \mid (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_{t-1}, y_{t-1}) \right] \end{aligned}$$

Applying Bernstein's gives

$$\frac{1}{n} \sum_{t=1}^n \text{risk}(H_{t-1}) \leq \frac{M_n}{n} + \frac{c}{n} \left(\ln M_n + \sqrt{M_n \ln M_n} \right) \quad \text{w.h.p.}$$

Recall

$$\frac{M_n}{n} = \frac{1}{n} \sum_{t=1}^n \mathbb{I}_{\{H_{t-1}(\mathbf{x}_t) \neq Y_t\}} \quad \text{is the fraction of mistakes}$$



Step 2: pick a good classifier in the ensemble

- Start from the ensemble H_0, H_1, \dots, H_n
- Do the following:
 - ① test each H_t on $(\mathbf{x}_{t+1}, \mathbf{y}_{t+1}), \dots, (\mathbf{x}_n, \mathbf{y}_n)$
 - ② pick $\hat{H} = H_{t^*}$ minimizing a **penalized risk estimate**

Guaranteed bound

$$\text{risk}(\hat{H}) \leq \frac{M_n}{n} + \frac{c}{n} \left((\ln n)^2 + \sqrt{M_n \ln n} \right) \quad \text{w.h.p.}$$



Conclusions!

- We made a reduction from pattern classification to prediction with experts
- Potential-based forecasters have on-line classifiers as natural counterparts
- We get a different viewpoint on kernel-based learning
- A simple large deviation inequality is enough to get data-dependent tail risk bounds

