

# Vovk's algorithm

## Mixable and unmixable loss functions

Yoav Freund

February 2, 2006

# Outline

Review

The general prediction game

Some useful loss functions

Vovk's algorithm

mixable loss functions

The convexity condition

Log loss

Square loss

    Square loss using simple averaging

Summary table

## The log-loss game

- ▶ Prediction algorithm  $A$  has access to  $N$  experts.
- ▶ The following is repeated for  $t = 1, \dots, T$ 
  - ▶ Experts generate predictive distributions:  $\mathbf{p}_1^t, \dots, \mathbf{p}_N^t$
  - ▶ Algorithm generates its own prediction  $\mathbf{p}_A^t$
  - ▶  $\mathbf{c}^t$  is revealed.
- ▶ **Goal:** minimize regret:

$$-\sum_{t=1}^T \log p_A^t(\mathbf{c}^t) + \min_{i=1, \dots, N} \left( -\sum_{t=1}^T \log p_i^t(\mathbf{c}^t) \right)$$

## The online Bayes Algorithm

- ▶ Total loss of expert  $i$

$$L_i^t = - \sum_{s=1}^t \log p_i^s(c^s); \quad L_i^0 = 0$$

- ▶ Weight of expert  $i$

$$w_i^t = w_i^1 e^{-L_i^{t-1}} = w_i^1 \prod_{s=1}^{t-1} p_i^s(c^s)$$

- ▶ Freedom to choose initial weights.

$$w_i^1 \geq 0, \sum_{i=1}^N w_i^1 = 1$$

- ▶ Prediction of algorithm  $A$

$$\mathbf{p}_A^t = \frac{\sum_{i=1}^N w_i^t \mathbf{p}_i^t}{\sum_{i=1}^N w_i^t}$$

## Cumulative loss vs. Final total weight

Total weight:  $W^t \doteq \sum_{i=1}^N w_i^t$

$$\frac{W^{t+1}}{W^t} = \frac{\sum_{i=1}^N w_i^t e^{\log p_i^t(c^t)}}{\sum_{i=1}^N w_i^t} = \frac{\sum_{i=1}^N w_i^t p_i^t(c^t)}{\sum_{i=1}^N w_i^t} = p_A^t(c^t)$$

$$-\log \frac{W^{t+1}}{W^t} = -\log p_A^t(c^t)$$

$$-\log W^{T+1} = -\log \frac{W^{T+1}}{W^1} = -\sum_{t=1}^T \log p_A^t(c^t) = L_A^T$$

**EQUALITY** not bound!

## Vovk's general prediction game

$\Gamma$  - prediction space.  $\Omega$  - outcome space.

On each trial  $t = 1, 2, \dots$

1. Each expert  $i \in \{1 \dots N\}$  makes a prediction  $\gamma_i^t \in \Gamma$
2. The learner, after observing  $\langle \gamma_1^t \dots \gamma_N^t \rangle$ , makes its own prediction  $\gamma^t$
3. Nature chooses an outcome  $\omega^t \in \Omega$
4. Each expert incurs loss  $\ell_i^t = \lambda(\omega^t, \gamma_i^t)$   
The learner incurs loss  $\ell_A^t = \lambda(\omega^t, \gamma^t)$

## Achievable loss bounds

- ▶  $L_A \doteq \sum_{t=1}^T \ell_A^t$  - total loss of algorithm
- ▶  $L_i \doteq \sum_{t=1}^T \ell_i^t$  - total loss of expert  $i$
- ▶ **Goal:** find an algorithm which guarantees that

$$(a, c) \in [0, \infty), \quad L_A \leq aL_{\min} + c \ln N$$

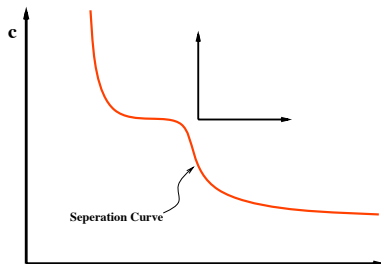
For any sequence of events.

- ▶ We say that the pair  $(a, c)$  is **achievable**.

## The set of achievable bounds

- ▶ Fix loss function  $\lambda : \Omega \times \Gamma \rightarrow [0, \infty)$
- ▶ The pair  $(a, c)$  is *achievable* if there exists *some* prediction algorithm such that for *any*  $N > 0$ , *any* set of  $N$  prediction sequences and *any* sequence of outcomes

$$L_A \leq aL_{\min} + c \ln N$$





## Some useful loss functions

- Outcomes:  $\omega^1, \omega_2, \dots \omega^t \in [0, 1]$
- Predictions:  $\gamma^1, \gamma^2, \dots \gamma^t \in [0, 1]$

## Log loss (Entropy loss)



$$\lambda_{\text{ent}}(\omega, \gamma) = \omega \ln \frac{\omega}{\gamma} + (1 - \omega) \ln \frac{1 - \omega}{1 - \gamma}$$

- ▶ When  $q_t \in \{0, 1\}$  Cumulative log loss = coding length  $\pm 1$
- ▶ If  $P[\omega_t = 1] = q$ , optimal prediction  $\gamma^t = q$
- ▶ Unbounded loss.
- ▶ Not symmetric  $\exists p, q \lambda(p, q) \neq \lambda(q, p)$ .
- ▶ No triangle inequality  
 $\exists p_1, p_2, p_3 \lambda(p_1, p_3) > \lambda(p_1, p_2) + \lambda(p_2, p_3)$

## Square loss (Breier Loss)



$$\lambda_{\text{sq}}(\omega, \gamma) = (\omega - \gamma)^2$$

- ▶  $P[\omega^t = 1] = q$ ,  $P[\omega^t = 0] = 1 - q$ ,  
optimal prediction  $\gamma^t = q$
- ▶ Bounded loss.
- ▶ Defines a metric (symmetric and triangle ineq.)
- ▶ Corresponds to regression.

## Hellinger Loss



$$\lambda_{\text{hel}}(\omega, \gamma) = \frac{1}{2} \left( (\sqrt{\omega} + \sqrt{\gamma})^2 + (\sqrt{1-\omega} + \sqrt{1-\gamma})^2 \right)$$

- ▶ If  $P[\omega^t = 1] = q$ ,  $P[\omega^t = 0] = 1 - q$ ,  
optimal prediction  $\gamma^t = q$
- ▶ Loss is bounded.
- ▶ Defines a metric.
- ▶  $\lambda_{\text{hel}}(p, q) \approx \lambda_{\text{ent}}(p, q)$  when  $p \approx q$  and  $p, q \in (0, 1)$

## Absolute loss



$$\lambda(\omega, \gamma) = |\omega - \gamma|$$

- ▶ Probability of making a mistake if predicting 0 or 1 using a biased coin
- ▶ If  $P[\omega^t = 1] = q$ ,  $P[\omega^t = 0] = 1 - q$ , then the optimal prediction is

$$\gamma^t = \begin{cases} 1 & \text{if } q > 1/2, \\ 0 & \text{otherwise} \end{cases}$$

## Structureless bounded loss

- ▶ Prediction is a distribution  $\gamma = \langle p_1, \dots, p_N \rangle$ ,  $p_i \geq 0$ ,  
 $\sum_{i=1}^N p_i = 1$
- ▶ Outcome is a loss vector  $\omega = \langle \omega_1, \dots, \omega_N \rangle$ ,  $0 \leq \omega_i \leq 1$
- ▶ Loss is the dot product:  $\lambda_{\text{dot}}(\omega, \gamma) = \gamma \cdot \omega$
- ▶ Corresponds to the hedging game.
- ▶ For hedge loss the regret is  $\Omega(\sqrt{T \log N})$ .
- ▶ For the log loss the regret is  $O(\log N)$
- ▶ Which losses behave like **entropy loss** and which behave like **hedge loss**?

## Some technical requirements

- ▶ There should be a **topology** on the prediction set  $\Gamma$  such that
- ▶  $\Gamma$  is compact.
- ▶  $\forall \omega \in \Omega$ , the function  $\gamma \rightarrow \lambda(\omega, \gamma)$  is **continuous**
- ▶ There is a **universally reasonable prediction**  
 $\exists \gamma \in \Gamma, \forall \omega \in \Omega, \lambda(\omega, \gamma) < \infty$
- ▶ There is **no universally optimal prediction**  
 $\neg \exists \gamma \in \Gamma, \forall \omega \in \Omega, \lambda(\omega, \gamma) = 0$

## Vovk's meta-algorithm

- Fix an **achievable** pair  $(a, c)$  and set  $\eta = a/c$

- 

- 1.

$$W_i^t = \frac{1}{N} e^{-\eta L_i^t}$$

2. Choose  $\gamma_t$  so that, for all  $\omega^t \in \Omega$ :

$$\lambda(\omega^t, \gamma^t) - c \ln \sum_i W_i^t \leq -c \ln \left( \sum_i W_i^t e^{-\eta \lambda(\omega^t, \gamma_i^t)} \right)$$

- If choice of  $\gamma^t$  always exists, then the total loss satisfies:

$$\sum_t \lambda(\omega^t, \gamma^t) \leq -c \ln \sum_i W_i^{T+1} \leq a L_{\min} + c \ln N$$

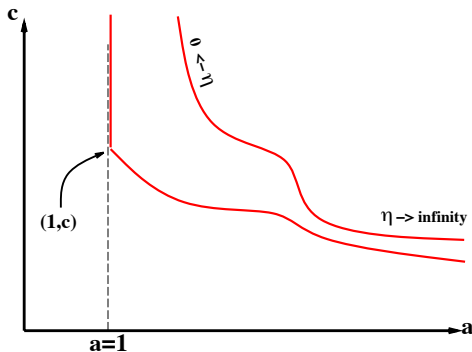
- Vovk's result: **yes!** a good choice for  $\gamma_t$  always exists!



# Vovk's algorithm is the the highest achiever [Vovk95]

The pair  $(a, c)$  is achieved by some algorithm if and only if it is achieved by Vovk's algorithm.

The separation curve is  $\left\{ \left( a(\eta), \frac{a(\eta)}{\eta} \right) \mid \eta \in [0, \infty] \right\}$



## Mixable Loss Functions

- ▶ A Loss function is **mixable** if a pair of the form  $(1, c)$ ,  $c < \infty$  is achievable.

$$L_A \leq L_{\min} + c \ln N$$

- ▶ Vovk's algorithm with  $\eta = 1/c$  achieves this bound.
- ▶  $\lambda_{\text{ent}}, \lambda_{\text{sq}}, \lambda_{\text{hel}}$  are **mixable**
- ▶  $\lambda_{\text{abs}}, \lambda_{\text{dot}}$  are **not mixable**

## The convexity condition

- ▶ requirement for loss to be  $(1, 1/\eta)$  mixable
- ▶  $\forall \langle (\gamma_1, W_1), \dots, (\gamma_N, W_N) \rangle$   
 $\exists \gamma \in \Gamma$   
 $\forall \omega \in \Omega$ :

$$\lambda(\omega, \gamma) - \frac{1}{\eta} \ln \sum_i W_i \leq -\frac{1}{\eta} \ln \left( \sum_i W_i e^{-\eta \lambda(\omega, \gamma_i)} \right)$$

- ▶ Can be re-written as:

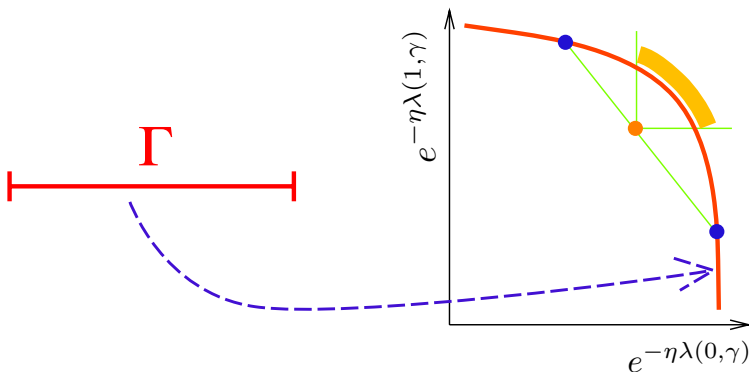
$$e^{-\eta \lambda(\omega, \gamma)} \geq \sum_i \left( \frac{W_i}{\sum_j W_j} \right) e^{-\eta \lambda(\omega, \gamma_i)}$$

- ▶ Equivalently - the image of the set  $\Gamma$  under the mapping  $F(\gamma) = \langle e^{-\eta \lambda(\omega, \gamma)} \rangle_{\omega \in \Omega}$  is concave.

## convexity condition: Pictorially

- **Example:** Suppose  $\Omega = \{0, 1\}$ ,  $\Gamma = [0, 1]$ . then

$$F(\gamma) = \langle e^{-\eta\lambda(0,\gamma)}, e^{-\eta\lambda(1,\gamma)} \rangle$$



## Vovk Algorithm for log loss

- ▶ The log loss is mixable with  $\eta = 1$
- ▶ The image of  $[0, 1]$  through  $F(\gamma) = \langle e^{-\eta\lambda(0,\gamma)}, e^{-\eta\lambda(1,\gamma)} \rangle$  is a straight line segment.
- ▶ The **only** satisfactory prediction is

$$\gamma = \frac{\sum_i w_i \gamma_i}{\sum_i w_i}$$

- ▶ We are back to the online Bayes algorithm.

## Vovk algorithm for square loss

- ▶ The square loss is mixable with  $\eta = 2$ .
- ▶ Prediction must satisfy

$$1 - \sqrt{-\frac{1}{2} \ln \sum_i V_i^t e^{-2(1-p_i^t)^2}} \leq p^t \leq \sqrt{-\frac{1}{2} \ln \sum_i V_i^t e^{-2(p_i^t)^2}}$$

where  $V_i^t = \frac{W_i^t}{\sum_s W_i^s}$ .



$$L_A \leq L_{\min} + \frac{1}{2} \ln N$$

## Simple prediction for square loss

- ▶ We can use the prediction

$$\gamma = \frac{\sum_i W_i \gamma_i}{\sum_i W_i}$$

- ▶ But in that case we must use  $\eta = 1/2$  when updating the weights.
- ▶ Which yields the bound

$$L_A \leq L_{\min} + 2 \ln N$$

## Summary of bounds for mixable losses

### TRACKING THE BEST EXPERT

Loss Functions:	$c$ values: $(\eta = 1/c)$	
	$\text{pred}_{\text{wmean}}(v, x)$	$\text{pred}_{\text{Vovk}}(v, x)$
$L_{\text{sq}}(p, q)$	2	$1/2$
$L_{\text{ent}}(p, q)$	1	1
$L_{\text{hel}}(p, q)$	1	$1/\sqrt{2}$

Figure 2.  $(c, 1/c)$ -realizability:  $c$  values for loss and prediction function pairings