

Tutorial on Universal Prediction, Day 2, morning



Peter Grünwald



Centrum Wiskunde & Informatica – Amsterdam
Mathematisch Instituut – Universiteit Leiden

Menu - Yesterday

1. Universal Prediction
with log-loss scoring
2. Universal Prediction and Bayesian Inference

Menu – This Morning

1. From log-loss to 0/1 loss
2. Bayesian Prediction, Precisely defined
3. Two Problems with Bayes Prediction for 0/1 loss
4. Generalizing Bayes, making it work for general loss
5. (Luckiness/Relevance)



Universal Prediction with log-loss



- On each i (day), Marjon and Peter **announce the probability** that $y_{i+1} = 1$, i.e. that it will rain on day $i+1$
- We would like to combine their predictions in some way such that for **every** sequence $y_1, \dots, y_n \in \{0, 1\}^n$ we predict almost as well as whoever turns out to be the best forecaster for that sequence in terms of their cumulative log-loss
 - If, with hindsight, Marjon was better, we predict as well as Marjon
 - If, with hindsight, Peter was better, we predict as well as Peter



Universal Prediction with 0/1-loss



- On each i (day), Marjon and Peter predict whether it will rain on day $i+1$, i.e. **they announce '1' or '0'**
- If their prediction is wrong, their loss is 1, otherwise 0



Universal Prediction with 0/1-loss



- On each i (day), Marjon and Peter predict whether it will rain on day $i+1$, i.e. **they announce '1' or '0'**
- If their prediction is wrong, their loss is 1, otherwise 0
- We would like to combine their predictions in some way such that for **every** sequence $y_1, \dots, y_n \in \{0, 1\}^n$ we predict almost as well as whoever turns out to be the best forecaster for that sequence in terms of cumulative 0/1-loss (**total nr of mistakes**)
 - If, with hindsight, Marjon was better, we predict as well as Marjon
 - If, with hindsight, Peter was better, we predict as well as Peter

Universal prediction with log loss

- We would like to combine predictions such that for every sequence $y_1, \dots, y_n \in \{0, 1\}^n$ we predict almost as well as the best forecaster for that sequence
- It turns out that there exists a universal strategy \bar{S} such that, for all $n, y_1, \dots, y_n \in \{0, 1\}^n$

$$\text{loss}(y_1 \dots, y_n, \bar{S}) \leq \min\{\text{loss}(y_1 \dots, y_n, S_{\text{Marjon}}), \text{loss}(y_1 \dots, y_n, S_{\text{Peter}})\} + 1.$$

Universal prediction with 0/1-loss

- We would like to combine predictions such that for every sequence $y_1, \dots, y_n \in \{0, 1\}^n$ we predict almost as well as the best forecaster for that sequence
- It turns out that there exists a universal strategy \bar{S} such that, for all $n, y_1, \dots, y_n \in \{0, 1\}^n$

$$\text{loss}(y_1 \dots, y_n, \bar{S}) \leq \min\{\text{loss}(y_1 \dots, y_n, S_{\text{Marjon}}), \text{loss}(y_1 \dots, y_n, S_{\text{Peter}})\} + \frac{X}{\sqrt{n}}$$

Beyond Bayes

- Good Algorithm for log-loss was Bayesian
 - Algorithm is minimax optimal up to constant factor (Vovk '99)
- Bayes can fail dramatically for 0/1-loss (and somewhat less dramatically for many other loss functions)
- Yet, intriguingly, a "simple" modification of Bayes is again essentially minimax optimal for general loss functions, including 0/1

Menu – This Morning

1. From log-loss to 0/1 loss
2. Bayesian Prediction, Precisely defined
3. Two Problems with Bayes Prediction for 0/1 loss
4. Generalizing Bayes and making it work for general losses
5. Final Remarks about General Losses

Bayesian Probability

- Let $\mathcal{M} = \{P_\theta : \theta \in \Theta\}$ be a finite or countable set of prob. distributions on $\mathcal{Y}^{\mathbb{N}}$.
- Let W be a distribution on Θ
- Suppose we believe that $\bar{\theta} \sim W$ and that $Y_1, Y_2, \dots \sim P_{\bar{\theta}}$

Bayesian Probability

- Let $\mathcal{M} = \{P_\theta : \theta \in \Theta\}$ be a finite or countable set of prob. distributions on $\mathcal{Y}^{\mathbb{N}}$.
- Let W be a distribution on Θ
- Suppose we believe that $\bar{\theta} \sim W$ and that $Y_1, Y_2, \dots \sim P_{\bar{\theta}}$
- This defines distribution P on joint space $\Theta \times \mathcal{Y}^{\mathbb{N}}$ with $P(\theta) := W(\theta) ; P(y^n | \theta) := P_\theta(y^n) ; P(\theta; y^n) := P(\theta) \cdot P(y^n | \theta)$

Bayesian Probability

- Let $\mathcal{A} = \{P_\theta : \theta \in \Theta\}$ be a **finite** or **countable** set of prob. distributions on \mathcal{Y}^{∞} .
- Let W be a distribution on Θ
- Suppose we believe that $\bar{\theta} \sim W$ and that $Y_1, Y_2, \dots \sim P_{\bar{\theta}}$
- This defines distribution P on joint space $\Theta \times \mathcal{Y}^{\infty}$ with

$$P(y^n) = \sum_{\theta \in \Theta} P(\theta)P(y^n | \theta) = \sum_{\theta \in \Theta} W(\theta)P_\theta(y^n) \leftarrow \text{marginal}$$

$$P(y_i | y^{i-1}) = P(y^i) / P(y^{i-1}) \leftarrow \text{predictive}$$

$$P(\theta | y^n) = P(y^n | \theta) \cdot P(\theta) / P(y^n) \leftarrow \text{posterior}$$

The Predictive and The Posterior

- The predictive distribution can be rewritten as

$$P(y_{i+1} | y_1, \dots, y_i) = \frac{P(y^{i+1})}{P(y^i)} = \frac{\sum_{\theta \in \Theta} P_\theta(y_{i+1} | y^i) P_\theta(y^i) W(\theta)}{\sum_{\theta \in \Theta} P_\theta(y^i) W(\theta)}$$

$$= \sum_{\theta \in \Theta} P_\theta(y_{i+1} | y^i) \cdot \frac{P_\theta(y^i) W(\theta)}{\sum_{\theta' \in \Theta} P_{\theta'}(y^i) W(\theta')} = \sum_{\theta \in \Theta} P_\theta(y_{i+1} | y^i) \cdot W(\theta | y^i)$$

Bayesian Prediction for general loss functions

- Let **loss** : $\mathcal{Y} \times \mathcal{A} \rightarrow [0, \infty]$ be arbitrary loss fn.
 - Log-loss: $\mathcal{A} = \Delta(\mathcal{Y}), \text{loss}(y, p) = -\log p(y)$
 - 0/1-loss: $\mathcal{A} = \mathcal{Y} = \{0, 1\}, \text{loss}(y, a) = |y - a|$

Bayesian Prediction for general loss functions

- Let **loss** : $\mathcal{Y} \times \mathcal{A} \rightarrow [0, \infty]$ be arbitrary loss fn.
 - Log-loss: $\mathcal{A} = \Delta(\mathcal{Y}), \text{loss}(y, p) = -\log p(y)$
 - 0/1-loss: $\mathcal{A} = \mathcal{Y} = \{0, 1\}, \text{loss}(y, a) = |y - a|$
- The Bayes prediction based on data y^i is given by

$$a_{\text{Bayes}} := \arg \min_{a \in \mathcal{A}} \mathbf{E}_{Y_{i+1} \sim P|y^i} [\text{loss}(Y_{i+1}, a)]$$

Bayesian Prediction for general loss functions

- Let **loss** : $\mathcal{Y} \times \mathcal{A} \rightarrow [0, \infty]$ be arbitrary loss fn.
 - Log-loss: $\mathcal{A} = \Delta(\mathcal{Y}), \text{loss}(y, p) = -\log p(y)$
- The Bayes prediction based on data y^i is given by

$$a_{\text{Bayes}} := \arg \min_{a \in \mathcal{A}} \mathbf{E}_{Y_{i+1} \sim P|y^i} [\text{loss}(Y_{i+1}, a)]$$

Log loss is a **proper scoring rule**: $p_{\text{Bayes}} = p(\cdot | y^{i-1})$

Proper Scoring Rule: if you believe P , say P !

Bayesian Prediction for general loss functions

- Let **loss** : $\mathcal{Y} \times \mathcal{A} \rightarrow [0, \infty]$ be arbitrary loss fn.
 - 0/1-loss: $\mathcal{A} = \mathcal{Y} = \{0, 1\}, \text{loss}(y, a) = |y - a|$
- The Bayes prediction based on data y^i is given by

$$a_{\text{Bayes}} := \arg \min_{a \in \mathcal{A}} \mathbf{E}_{Y_{i+1} \sim P|y^i} [\text{loss}(Y_{i+1}, a)]$$

0/1 loss:

Bayesian Prediction for general loss functions

- Let **loss** : $\mathcal{Y} \times \mathcal{A} \rightarrow [0, \infty]$ be arbitrary loss fn.
- 0/1-loss: $\mathcal{A} = \mathcal{Y} = \{0, 1\}$, $\text{loss}(y, a) = |y - a|$
- The Bayes prediction based on data y^i is given by

$$a_{\text{Bayes}} := \arg \min_{a \in \mathcal{A}} \mathbf{E}_{Y_{i+1} \sim P|y^i} [\text{loss}(Y_{i+1}, a)]$$

$$\text{0/1 loss: } a_{\text{Bayes}} = \begin{cases} 0 & \text{if } P(y_i = 1 | y^{i-1}) < \frac{1}{2} \\ 1 & \text{if } P(y_i = 1 | y^{i-1}) > \frac{1}{2} \end{cases}$$

Bayes vs. Worst-Case

- With **log-loss**, the algorithm presented yesterday really behaves identical to Bayesian sequential prediction and we get great $-\log W(\theta)$ regret bound
- With **0/1-loss**, there are **2 cases...** and **2 problems:**
 - Predictors issue probabilities (yesterday's Peter & Margot), to be used for 0/1-predictions \rightarrow Bayes can be applied, **but now falls into similar traps as follow the leader!**
 - Predictors directly issue 0s and 1s (today's Peter & Margot) \rightarrow **Bayes cannot even be applied**

Menu – This Morning

- From log-loss to 0/1 loss
- Bayesian Prediction, Precisely defined
- Two Problems with Bayes Prediction for 0/1 loss**
- Generalizing Bayes, making it work for general loss
- Additional Remarks

Recall the FTL Problem

- FTL**: at each day n ,
 - if Marjon was best on y_1, \dots, y_n , then predict y_n , | 1 like Marjon
 - if Peter was best on y_1, \dots, y_n then predict y_n , | 1 like Peter
- Proposition: there exist Peter and Marjon such that

$$\max_{y_1, \dots, y_n} \{ \text{loss}(y_1, \dots, y_n, \bar{S}) - \min_{S \in \text{Peter, Marjon}} \text{loss}(y_1, \dots, y_n, S) \} = 0.25n$$

Each day Peter says 'it rains with probability 1/4'
 Marjon says 'it rains with probability 3/4'

010101010101010

Bayes' Problem

- After every odd day, posterior probability 'on' Peter will be **slightly larger than 1/2**: $P(y_i | y^{i-1}) = 0.5 + \tilde{O}(1/i)$
- Then Bayes will predict 0...so Bayes predicts like FTL**

Each day Peter says 'it rains with probability 1/4'
 Marjon says 'it rains with probability 3/4'

010101010101010

Bayes vs. Worst-Case

- With **log-loss**, the algorithm presented yesterday really behaves identical to Bayesian sequential prediction and we get great $-\log W(\theta)$ regret bound
- With **0/1-loss**, there are **2 cases...** and **2 problems:**
 - Predictors issue probabilities (yesterday's Peter & Margot), to be used for 0/1-predictions \rightarrow Bayes can be applied, but now falls into similar traps as follow the leader!
 - Predictors directly issue 0s and 1s (today's Peter & Margot) \rightarrow **Bayes cannot even be applied**

Bayes vs. Worst-Case

- With **log-loss**, the algorithm presented yesterday really behaves identical to Bayesian sequential prediction and we get great $-\log W(\theta)$ regret bound
- With **0/1-loss**, there are **2 cases...** and **2 problems**:
 1. Predictors issue probabilities (yesterday's Peter & Margot), to be used for 0/1-predictions \rightarrow Bayes can be applied, but now falls into similar traps as follow the leader!
 2. Predictors directly issue 0s and 1s (today's Peter & Margot) \rightarrow **Bayes cannot even be applied**

Strategy: first attack problem 2, then get back to 1
 Result is generalization of Bayes that is close to minimax optimal for arbitrary loss fns, but regret bounds do get worse

Solving Problem 2: Generalizing Bayes

- The standard solution consists of a mapping from arbitrary loss functions to the log loss that I called **entropification** in earlier work (G., COLT 1998)

Menu – This Morning

1. From log-loss to 0/1 loss
2. Bayesian Prediction, Precisely defined
3. Two Problems with Bayes Prediction for 0/1 loss
4. **Generalizing Bayes**, making it work for general loss
5. Additional Remarks

Entropification (The Gauss Device)

- Let **loss** : $\mathcal{Y} \times \mathcal{A} \rightarrow [0, \infty]$ be arbitrary loss fn.
- For every action a and $\eta > 0$ define density

$$p_{a,\eta}(y) := \frac{1}{Z(\eta)} e^{-\eta \text{loss}(y,a)}$$

Entropification (The Gauss Device)

- Let **loss** : $\mathcal{Y} \times \mathcal{A} \rightarrow [0, \infty]$ be arbitrary loss fn.
- For every action a and $\eta > 0$ define density

$$p_{a,\eta}(y) := \frac{1}{Z(\eta)} e^{-\eta \text{loss}(y,a)} \quad Z(\eta, \mathbf{a}) = \int_{\mathcal{Y}} e^{-\eta \text{loss}(y,a)} dy$$

We can do this for just about every loss function, and will do it here for 0/1-loss. Complications with $Z(\eta)$ for nonsymmetric losses can be solved (G., 2008)

Entropification (The Gauss Device)

- Let **loss** : $\mathcal{Y} \times \mathcal{A} \rightarrow [0, \infty]$ be arbitrary loss fn.
- For every action a and $\eta > 0$ define density

$$p_{a,\eta}(y) := \frac{1}{Z(\eta)} e^{-\eta \text{loss}(y,a)} \quad Z(\eta, \mathbf{a}) = \int_{\mathcal{Y}} e^{-\eta \text{loss}(y,a)} dy$$

- Example: **0/1-loss**: $Z(\eta) = e^{-\eta \cdot 0} + e^{-\eta \cdot 1} = 1 + e^{-\eta}$

$$p_{a,\eta}(y) := \frac{1}{Z(\eta)} e^{-\eta \text{loss}(y,a)} = \begin{cases} \frac{e^{-\eta}}{1+e^{-\eta}} = \frac{1}{2} - b & \text{if } y \neq a \\ \frac{1}{1+e^{-\eta}} = \frac{1}{2} + b & \text{if } y = a \end{cases}$$

Entropification (The Gauss Device)

- Let **loss** : $\mathcal{Y} \times \mathcal{A} \rightarrow [0, \infty]$ be arbitrary loss fn.
- For every action a and $\eta > 0$ define density

$$p_{a,\eta}(y) := \frac{1}{Z(\eta)} e^{-\eta \text{loss}(y,a)} \quad Z(\eta, \mathbf{a}) = \int_{\mathcal{Y}} e^{-\eta \text{loss}(y,a)} dy$$

- Example: **squared loss**:

$$\text{loss}(y, a) = (y - a)^2, \mathcal{Y} = \mathcal{A} = \mathbb{R}$$

$$p_{a,\eta}(y) = \frac{1}{Z(\eta)} e^{-\eta(y-a)^2} = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{\sigma^2}(y-a)^2}$$

Entropification (The Gauss Device)

- log loss of constructed distributions is affine (linear+constant) function of loss of interest

$$\text{log-loss}(y, p_{a,\eta}) = -\log p_{a,\eta}(y) = \eta \text{loss}(y, S) + \ln Z(\eta)$$

- log-loss difference is even linear in loss-difference:

$$\text{log-loss}(y, p_{a,\eta}) - \text{log-loss}(y, p_{a',\eta}) = \eta (\text{loss}(y, a) - \text{loss}(y, a'))$$

action is good for original loss iff transformed action is good for log-loss!

Entropification (The Gauss Device)

- Let **loss** : $\mathcal{Y} \times \mathcal{A} \rightarrow [0, \infty]$ be arbitrary loss fn.
- For every action a and $\eta > 0$ define prob. mass fn.

$$p_{a,\eta}(y) := \frac{1}{Z(\eta)} e^{-\eta \text{loss}(y,a)}$$

- A **strategy** relative to \mathcal{A} is a function $S : (\mathcal{X} \times \mathcal{Y})^* \rightarrow \mathcal{A}$

Entropification (The Gauss Device)

- Let **loss** : $\mathcal{Y} \times \mathcal{A} \rightarrow [0, \infty]$ be arbitrary loss fn.
- For every action a and $\eta > 0$ define prob. mass fn.

$$p_{a,\eta}(y) := \frac{1}{Z(\eta)} e^{-\eta \text{loss}(y,a)}$$

- A **strategy** relative to \mathcal{A} is a function $S : (\mathcal{X} \times \mathcal{Y})^* \rightarrow \mathcal{A}$
- Extend definition to strategies as:

$$p_{S,\eta}(y_i | x^{i-1}, y^{i-1}) := p_{S(x^{i-1}, y^{i-1}), \eta}(y_i)$$

$$p_{S,\eta}(y^n | x^n) := \prod_{i=1}^n p_{S,\eta}(y_i | y^{i-1}, x^{i-1})$$

$$= \frac{1}{Z(\eta)^n} \cdot e^{-\eta \sum_{i=1}^n \text{loss}(y_i, S(y^{i-1}, x^{i-1}))} = \frac{1}{Z(\eta)^n} e^{-\eta \text{loss}(y^n, S)}$$

Entropification (The Gauss Device)

- Let **loss** : $\mathcal{Y} \times \mathcal{A} \rightarrow [0, \infty]$ be arbitrary loss fn.
- For every action a and $\eta > 0$ define prob. mass fn.

$$p_{a,\eta}(y) := \frac{1}{Z(\eta)} e^{-\eta \text{loss}(y,a)}$$

- A **strategy** relative to \mathcal{A} is a function $S : (\mathcal{X} \times \mathcal{Y})^* \rightarrow \mathcal{A}$
- Extend definition to strategies as:

$$p_{S,\eta}(y^n | x^n) = \frac{1}{Z(\eta)^n} \cdot e^{-\eta \text{loss}(y^n, S)}$$

Entropification (The Gauss Device)

- accumulated log loss prediction error is affine (linear+constant) function of loss of interest

$$\text{log-loss}(y^n, p_{S,\eta}) = -\log p_{S,\eta}(y^n) = \eta \text{loss}(y^n, S) + n \ln Z(\eta)$$

- accumulated loss **difference** is even linear:

$$\text{log-loss}(y^n, p_{S,\eta}) - \text{log-loss}(y^n, p_{S',\eta}) = \eta (\text{loss}(y^n, S) - \text{loss}(y^n, S'))$$

prediction strategy is good for original loss iff transformed prediction strategy is good for log-loss!

Applying Bayes to general predictors

- Recall:

$$p_{\theta, \eta}(y_i | y^{i-1}) = \frac{1}{Z(\eta)} e^{-\eta \text{loss}(y_i, \theta(y^{i-1}))}$$
- We now get

$$P_{\text{Bayes}, \eta}(y_{i+1} | y^i) = \sum_{\theta} W_{\eta}(\theta | y^i) \cdot \left(\frac{1}{Z(\eta)} e^{-\eta \text{loss}(y_i, \theta(y^{i-1}))} \right)$$
- ... with the **generalized posterior**

$$W_{\eta}(\theta | y^i) = \frac{e^{-\eta \text{loss}(y^i, \theta)} \cdot W(\theta)}{\sum_{\theta' \in \Theta} e^{-\eta \text{loss}(y^i, \theta')} \cdot W(\theta')}$$

With $\eta = 1$ and log-loss this reduces to standard Bayes posterior

Bayes is good with log-loss

- For all n, y^n , all ℓ :

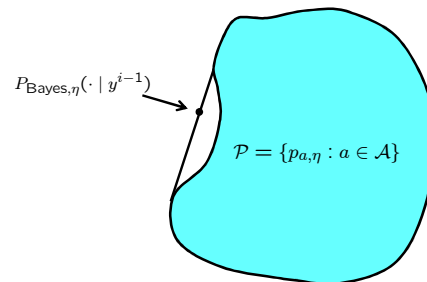
$$\log\text{-loss}(y^n, P_{\text{Bayes}}) \leq \log\text{-loss}(y^n, P_{\theta}) - \log W(\theta)$$
- For all sequences of each length n , **regret** of Bayes bounded by constant depending on ℓ , not on n
- For "nonmixable" loss functions like 0/1-loss and absolute loss, this does not work
 - standard Bayes does not perform well at all in worst-case
 - Optimal algorithm gets much larger regret of order $\sqrt{n(-\log W(\theta))}$ in worst-case

WHY??

Answer: Bayes goes beyond model

- Bayesian predictive distribution steps outside "model" $\mathcal{P} = \{p_{a, \eta} : a \in \mathcal{A}\}$: it predicts by **mixture** of $p_{a, \eta}$
- But if we want a prediction strategy applicable for original loss, we must always predict by p of form $p_{a, \eta} \propto \exp(-\eta \text{loss})$

Bayes goes Beyond Model



Forcing Bayes into the Model

- Bayesian predictive distribution steps outside model : it predicts by **mixture** of $p_{a, \eta}$
- But if we want a prediction strategy applicable for original loss, we must always predict by p of form $p_{a, \eta} \propto \exp(-\eta \text{loss})$
- We then need some algorithm **A** to turn Bayes posterior into allowed prediction. Examples:
 - Predict by **MAP**:

$$\log\text{-loss}(y_i, \text{map}(W(\cdot | y^{i-1}))) := -\log p_{\theta_{\text{map}}}(y_i | \theta)$$
 where θ_{map} achieves maximum of $W(\theta | y^{i-1})$
 - Predict by $p_{a, \eta}$ minimizing **posterior expected loss**

The Most Important Notion: $\Delta_{\eta}^*(\mathbf{A})$

- We define the **mixability gap** (G. et al. 2011) to be:

$$\Delta_{\eta}^*(\mathbf{A}) = \log\text{loss}(y^n, \mathbf{A}) - \log\text{loss}(y^n, P_{\text{Bayes}, \eta})$$
- Interpretation: amount of bits (loss units) lost by being forced to use allowed predictions instead of using the happily mixing Bayes prediction

The Clue

- We then have for all n, y^n, θ :

$$\log\text{-loss}(y^n, P_{\text{Bayes}, \eta}) \leq \log\text{-loss}(y^n, P_{\theta, \eta}) - \log W(\theta)$$

$$\log\text{-loss}(y^n, \mathbf{A}) \leq \log\text{-loss}(y^n, \theta) - \log W(\theta) + \Delta_{\eta}^*(\mathbf{A})$$

$$\text{loss}(y^n, \mathbf{A}) \leq \text{loss}(y^n, \theta) + \frac{-\log W(\theta)}{\eta} + \frac{\Delta_{\eta}^*(\mathbf{A})}{\eta}$$

The Clue

- We then have for all n, y^n, θ :

$$\log\text{-loss}(y^n, P_{\text{Bayes}, \eta}) \leq \log\text{-loss}(y^n, P_{\theta, \eta}) - \log W(\theta)$$

$$\log\text{-loss}(y^n, \mathbf{A}) \leq \log\text{-loss}(y^n, \theta) - \log W(\theta) + \Delta_{\eta}^*(\mathbf{A})$$

$$\text{loss}(y^n, \mathbf{A}) \leq \text{loss}(y^n, \theta) + \frac{-\log W(\theta)}{\eta} + \frac{\Delta_{\eta}^*(\mathbf{A})}{\eta}$$

Vovk '90: for so-called "mixable" loss fns, there exists an \mathbf{A} (the aggregating algorithm) such that for some $\eta > 0$, we are guaranteed $\Delta_{\eta}^*(\mathbf{A}) \leq 0$ (hence name **mixability gap**)

The Clue

- We then have for all n, y^n, θ :

$$\log\text{-loss}(y^n, P_{\text{Bayes}, \eta}) \leq \log\text{-loss}(y^n, P_{\theta, \eta}) - \log W(\theta)$$

$$\log\text{-loss}(y^n, \mathbf{A}) \leq \log\text{-loss}(y^n, \theta) - \log W(\theta) + \Delta_{\eta}^*(\mathbf{A})$$

$$\text{loss}(y^n, \mathbf{A}) \leq \text{loss}(y^n, \theta) + \frac{-\log W(\theta)}{\eta} + \frac{\Delta_{\eta}^*(\mathbf{A})}{\eta}$$

Vovk '90: for so-called "mixable" loss fns, there exists an \mathbf{A} (the aggregating algorithm) such that for some $\eta > 0$, we are guaranteed $\Delta_{\eta}^*(\mathbf{A}) \leq 0$ (hence name **mixability gap**)

Example: squared loss
 $\text{loss}(y, a) = (y - a)^2$ is $\frac{1}{2}$ -mixable if $\mathcal{Y} = [-1, 1]$

Mixable Loss Functions


- All 'strictly convex' loss functions with bounded range are mixable for finite $\eta > 0$
- So we can still run generalized Bayes and get regret bounds that are still of order $O(-\log W(\theta))$
- Usually $\eta < 1$:
prior becomes more, data less important
- Modifying Bayes may seem like 'hack' (relation to Bayes' theorem is lost) but: **resulting procedure still minimax optimal up to a constant (and in practice preferable over minimax algorithm)**
- Something deep going on...

The Clue, now for nonmixable loss

- We then have for all n, y^n, θ :

$$\text{loss}(y^n, \mathbf{A}) \leq \text{loss}(y^n, \theta) + \frac{-\log W(\theta)}{\eta} + \frac{\Delta_{\eta}^*(\mathbf{A})}{\eta}$$

0/1-loss is nonmixable
 (related to Bayes predictions being 'pushed to boundary')




The Clue

- We then have for all n, y^n, θ :

$$\text{loss}(y^n, \mathbf{A}) \leq \text{loss}(y^n, \theta) + \frac{-\log W(\theta)}{\eta} + \frac{\Delta_{\eta}^*(\mathbf{A})}{\eta}$$

Hedge-style algorithms: for bounded nonmixable losses functions: there still exist \mathbf{A} such that $\Delta_{\eta}^*(\mathbf{A}) \leq \eta^2 \cdot C$
 For uniform prior over K predictors, optimizing over η gives $\eta = \tilde{O}\left(\sqrt{\frac{\log K}{n}}\right)$
 We get regret bound (and actual regret) $O(\sqrt{n \cdot (\log K)})$
 (Warmuth, Freund, Schapire)



The Clue

- We then have for all n, y^n, θ :

$$\text{loss}(y^n, \mathbf{A}) \leq \text{loss}(y^n, \theta) + \frac{-\log W(\theta)}{\eta} + \frac{\Delta_{\eta}^*(\mathbf{A})}{\eta}$$

Hedge-style algorithms: for bounded nonmixable losses functions: there still exist \mathbf{A} such that $\Delta_{\eta}^*(\mathbf{A}) \leq \eta^2 \cdot C$
 For uniform prior over K predictors, optimizing over η gives $\eta = \mathcal{O}\left(\sqrt{\frac{\log K}{n}}\right)$
 We get regret bound (and actual regret) $\mathcal{O}\left(\sqrt{n \cdot (\log K)}\right)$
 (Warmuth, Freund, Schapire)

The Hedge Algorithm

- Basic Algorithm requires knowledge of 'horizon' n :
 at every time point $t < n$, you use generalized posterior with $\eta = \mathcal{O}\left(\sqrt{\frac{\log K}{n}}\right)$
- If n unknown, can still achieve same bound by decreasing learning rate dynamically: at each time t , you use posterior weights you get if you had used $\eta = \mathcal{O}\left(\sqrt{\frac{\log K}{t}}\right)$ from time 1 to t


"The older you get, the less attention you pay to all the experiences you had in life!"

Menu – This Morning

1. From log-loss to 0/1 loss
2. Bayesian Prediction, Precisely defined
3. Two Problems with Bayes Prediction for 0/1 loss
4. Generalizing Bayes, making it work for general loss
5. Final Remarks about 0/1-loss

Varying the Setting

- Predicting better than the Best Expert
 - As good as the best convex combination
 - As good as the best sequence of experts
 - Applied to Prediction of Electricity Consumption in Greater Paris Region by Electricité de France (Devaine, Goude, Stoltz 2012)
- Hedge with Infinitely Many Experts
- Bandit and Other Limited Feedback Settings



Major Open Problem(s)

Learning the Learning Rate

Major Open Problem(s)

Learning the Learning Rate

Easy Data