# Ensemble Learning - 2 Beyond Linear Aggregation

### John Klein



**Goal** : Resolve ties on super-majority / approval votes **Use case** : Recommender system aggregation



◆□▶ ◆□▶ ◆三▶ ◆三▶ ○○ ○○

### Instant Runoff Operator (IRO) :



#### Example

#### Fusion of recommender systems

- Suppose we trained *M* recommender systems each returning a list of #*C*=4 objects.
- Each list is ordered from most preferred to least preferred.
- Let us also assume that they delivered only 4 different kinds of ballots.
- The table in the next slide gives the advocacies along with their occurrences.
- Using the majority operator, site web A would win.

### Example

#### Ballot table for example

occurrences	42	26	15	17
ballot	Α	В	C	D
	В	C	D	C
	C	D	В	В
	D	Α	A	Α

▲□▶ ▲□▶ ▲ 臣▶ ★ 臣▶ 三臣 - のへぐ

Pairwise Rank Operator (PRO) :



▲ロト ▲圖ト ▲注ト ▲注ト 注 のへで

Pairwise Rank Operator (PRO) :

- ➤ "Weakest" win : find arg min V<sub>ij</sub> among those (i, j) corresponding to a win (W<sub>ij</sub> = 1).
- ► Let us run PRO on the recommender system fusion example.
- Neither PRO or IRO return a full list of preferred candidates. To obtain the list, eliminate the returned class from the ballots and re-run the algorithm.

#### Generalized averages :

 F-means are a class of generalized averaging obtained thanks to a bijective mapping g : Y → R.

$$f_{\mathrm{ens}}\left(\mathbf{x}
ight) = g^{-1}\left(rac{1}{M}\sum_{m=1}^{M}g\left(f_{m}\left(\mathbf{x}
ight)
ight)
ight).$$

- Typical choices are
  - g = log : geometric mean,
  - g = 1/x : harmonic mean.

◆ロト ◆昼 ▶ ◆臣 ▶ ◆臣 ● ○○○

Geometric mean :

$$f_{\mathsf{ens}}\left(\mathbf{x}\right) = \left(\prod_{m=1}^{M} f_{m}\left(\mathbf{x}\right)\right)^{\frac{1}{M}}$$

 Useful when aggregated values are normalized w.r.t some reference value (ratios or percentages).

#### Example

An orange tree produced in the past years 100, 180, 210 and 300 oranges. The corresponding growth production rates are thus 1.8, 1.17 and 1.44. The arithmetic mean of these rates is 1.47 whereas the geometric mean is 1.44. Starting with 100 oranges and applying a growth rate of 1.47 gives 314.47 oranges while using 1.44 as growth rate gives exactly 300 oranges.

Harmonic mean :

$$f_{\mathrm{ens}}\left(\mathbf{x}
ight) = rac{M}{\sum_{m=1}^{M}rac{1}{f_{m}(\mathbf{x})}}.$$

- Useful when aggregated values are unnormalized growth rates.
- ► Harmonic mean < Geometric mean < Arithmetic mean.

#### Example

A car is moving by a 100m at 50km/h and by another 100m at 70km/h. The car displacement duration is 7.2 for the first move followed by 5.14s for the second move which makes 12.34s in total.

- Arithmetic mean of (50; 70) is 60. The displacement duration obtained by moving at 60km/h on 200m is 12s.
- Geometric mean of (50; 70) is 59.16. The displacement duration obtained by moving at 59.16km/h on 200m is 12.17s.
- Harmonic mean of (50; 70) is 58.33. The displacement duration obtained by moving at 58.33km/h on 200m is exactly 12.34s.

Classification : combining predictions that are class label conditional probabilities

Aggregation of probability distributions : desirable properties

- ► Let *p*<sub>ens</sub> denote the aggregated distribution.
- ► Axiom (i) : weak set wise function property (WSFP)

Definition (WSFP)

for som

For all subset  $A \subseteq C$ ,

$$p_{ens}(A) = g^{(A)}(p_1(A), \dots, p_M(A)), \quad (1)$$
  
e function  $g^{(A)} : [0; 1]^M \to [0; 1].$ 

 Interpretation : the aggregated chances of event A are depending solely on the input probabilities on the same event A.

Classification : combining predictions that are class label conditional probabilities

Aggregation of probability distributions : desirable properties

► Axiom (ii) : strong set wise function property (SSFP)

## Definition (SSFP)

For all subset  $A \subseteq \mathcal{X}$ ,

$$p_{ens}(A) = g(p_1(A), \dots, p_M(A)),$$
 (2)

for some function  $g: [0; 1]^M \rightarrow [0; 1]$ .

Interpretation : same as before but the combination rule is the same for each event otherwise relabeling the elements of C would impact the fusion. Classification : combining predictions that are class label conditional probabilities

Aggregation of probability distributions : desirable properties

• Axiom (iii) : unanimity (or idempotence)

Definition (Unanimity) If  $p_m = p_0$  for all *m*, then  $p_{ens} = p_0$ .

Interpretation : if the sources are unanimous, then the aggregate distribution is a copy of the input ones.

### Proposition

If  $\#C \ge 3$ , a probability distribution aggregation operator satisfying SSFP and unanimity is an LOP (convex combination).

Classification : combining predictions that are class label conditional probabilities

Aggregation of probability distributions : desirable properties

Axiom (iv) : independence preservation (IP)

## Definition (IP)

For any two subsets A and B of C s.t.  $p_m(A \cap B) = p_m(A) \times p_m(B) \forall m$ , then  $p_{ens}(A \cap B) = p_{ens}(A) \times p_{ens}(B)$ 

- Interpretation : when input distributions are unanimously independent w.r.t. a pair of events, this is also true for the aggregated distribution.
- ► No LOP operator achieves IP except if w<sub>m</sub> = 1 for some m (dictatorship or selection).

Classification : combining predictions that are class label conditional probabilities

Aggregation of probability distributions : desirable properties

- In the Bayesian setting, the **posterior** p' is an update of the **prior** p through the **likelihood** function L : p'(x) ∝ L(D) p(x).
- ► Axiom (v) : Bayesian externality (EB)

## Definition (EB)

Let  $(p')_{ens}$  denote the combination of the updated distribution  $p'_m$  using likelihood function L and  $(p_{ens})'$  denote the updated combination of the distributions  $p_m$  using the same likelihood function. Then

$$(p')_{\mathsf{ens}} = (p_{\mathsf{ens}})'.$$

- Interpretation : Bayesian update and fusion commute.
- The time at which some information arrives does not matter.

Classification : combining predictions that are class label conditional probabilities

Aggregation of probability distributions : desirable properties

#### Proposition

If a probability distribution fusion operator writes

$$p_{ens}=rac{1}{Z}\prod_{m=1}^{M}\left(p_{m}
ight)^{w_{m}},$$

where coefficients  $w_m$  are non-negative and sum to one :  $\sum_{m=1}^{M} w_m = 1$  and Z is a normalization constant, then it achieves unanimity and EB.

- These operators are known as logarithmic opinion pool (LogOP) operators.
- LogOP = weighted geometric mean !

Classification : combining predictions that are class label conditional probabilities

Aggregation of probability distributions : desirable properties

• Another nice property of LogOPs :

### Proposition

An aggregated distribution returned by a LogOP is the solution of the minimization problem :

$$\arg\min_{p}\sum_{m=1}^{M}w_{m}d_{KL}\left(p,p_{m}\right),$$

(□) (@) (E) (E) E

where  $d_{KL}$  is the Kullback-Leibler divergence.

► LogOPs are "barycenters" in the KL divergence sense.

Classification : combining predictions that are class label conditional probabilities

Aggregation of probability distributions : probabilistic calculus

- Each distribution  $p_m$  can be interpreted as  $p_Y(.|\mathbf{x}, \mathbf{z}_m)$
- ►  $Z_m$  is a random variable attached to the  $m^{\text{th}}$  predictor, ex :  $Z_m = \mathcal{D}^{(m)}$ , private dataset of the  $m^{\text{th}}$  learner.
- For simplicity, we assume M = 2.
- From Bayesian standpoint, we want to infer  $p_Y(.|\mathbf{x}, \mathcal{D}^{(1)}, \mathcal{D}^{(2)}).$

• Let 
$$w_m = p\left(\mathcal{D}^{(m)}|\mathbf{x}\right)$$
.

▲ロト ▲圖ト ▲注ト ▲注ト 注 のへで

Classification : combining predictions that are class label conditional probabilities

Aggregation of probability distributions : probabilistic calculus

Applying Bayes, we can write

$$p\left(y|\mathbf{x}, \mathcal{D}^{(1)}, \mathcal{D}^{(2)}\right) \propto p\left(\mathcal{D}^{(1)}, \mathcal{D}^{(2)}|y, \mathbf{x}\right) \times p_0\left(y|\mathbf{x}\right)$$
(3)

Assuming conditional independence, we have

$$p\left(y|\mathbf{x}, \mathcal{D}^{(1)}, \mathcal{D}^{(2)}\right) \propto p\left(\mathcal{D}^{(1)}|y, \mathbf{x}\right) \times p\left(\mathcal{D}^{(2)}|y, \mathbf{x}\right) \times p_0\left(y|\mathbf{x}\right)$$
(4)

Applying Bayes again gives

$$p\left(y|\mathbf{x}, \mathcal{D}^{(1)}, \mathcal{D}^{(2)}\right) \propto \frac{p\left(y|\mathcal{D}^{(1)}, \mathbf{x}\right) w_1}{p_0\left(y|\mathbf{x}\right)} \times \frac{p\left(y|\mathcal{D}^{(2)}, \mathbf{x}\right) w_2}{p_0\left(y|\mathbf{x}\right)} \times p_0\left(y|\mathbf{x}\right)$$
(5)

Classification : combining predictions that are class label conditional probabilities

Aggregation of probability distributions : probabilistic calculus

This also writes

$$p\left(y|\mathbf{x}, \mathcal{D}^{(1)}, \mathcal{D}^{(2)}\right) \propto \frac{p\left(y|\mathcal{D}^{(1)}, \mathbf{x}\right) w_1 \times p\left(y|\mathcal{D}^{(2)}, \mathbf{x}\right) w_2}{p_0\left(y|\mathbf{x}\right)}$$
(6)

• 
$$p_0(y|\mathbf{x})$$
 is a prior  $\neq p(y|\mathbf{x})$ .

General case :

$$p\left(y|\mathbf{x}, \mathcal{D}^{(1)}, .., \mathcal{D}^{(M)}\right) \propto \frac{\prod_{m=1}^{M} p\left(y|\mathcal{D}^{(m)}, \mathbf{x}\right) w_{m}}{\left(p_{0}\left(y|\mathbf{x}\right)\right)^{M-1}}.$$
 (7)

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 ...

Classification : combining predictions that are class label conditional probabilities

Aggregation of probability distributions : when they are no longer defined on the same  $\sigma$ -algebra.

- ► Suppose f<sub>1</sub> is a trained logistic regressor to discriminate c<sub>1</sub> from B<sub>1</sub> = {c<sub>2</sub>, ..., c<sub>ℓ</sub>}
- The σ-algebra on which the distribution attached to f<sub>1</sub> is defined is : {∅, {c<sub>1</sub>}, B<sub>1</sub>, C}.
- ▶ If each *f<sub>m</sub>* is a one-versus-all probabilistic classifier, their distributions cannot be aggregated using the aforementioned techniques.
- Generalized frameworks can help : imprecise probabilities, belief functions/random sets.

# Learning to Aggregate Classification





## Learning to Aggregate Classification

**Stacking** : Procedure Generate a new training set for a second-stage classifier.



<□ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

# Learning to Aggregate Classification

### **Stacking** : Procedure Train the second-stage classifier.



< ロ > < 同 > < 三 > < 三 > < 三 > < ○ < ○ </p>

#### Learning to Aggregate Classification

# Stacking : Procedure Use it all.



990

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

### Learning to Aggregate Classification

### Stacking : Comments

- ► Stacking can be used for heterogeneous classifiers too.
- ► Need for validation data for second-stage training.
- Solution : repeat the procedure in a cross validation (CV) fashion.
- Leave-one-out CV (LOOCV) : let f<sub>m</sub><sup>(-i)</sup> denote the m<sup>th</sup> predictor trained on the whole training set except data point x<sup>(i)</sup>, then in step 3, one needs to minimize

$$J(\mathbf{w}) = \sum_{i=1}^{n} L\left(y^{(i)}, \operatorname{sgm}\left(\sum_{m=1}^{M} w_m f_m^{(-i)}\left(\mathbf{x}^{(i)}\right)\right)\right).$$

# Learning to Aggregate Classification

Any optimal solution?

• We assume  $f_m(\mathbf{x}) \in C$  (predictions are class labels).

• Let 
$$\mathbf{c} = [f_1(\mathbf{x}) \dots f_M(\mathbf{x})].$$

 Based on the sole information contained in vector c, and under 0-1 loss, the optimal decision is

$$f_{\text{ens}}(\mathbf{x}) = \underset{y}{\arg \max} p(y|\mathbf{c}(\mathbf{x})).$$

- Can we infer the conditional distributions  $p(y|\mathbf{c}(\mathbf{x}))$ ?
- Applying Bayes, we have  $p(y|\mathbf{c}) \propto p(\mathbf{c}|y) p(y)$ .
- $p(\mathbf{c}|y)$  is a joint distribution whose marginals are multinomial.
- How many parameters do we have to learn?

# Learning to Aggregate Classification

Any optimal solution?

- ► Not scalable : polynomial memory complexity w.r.t. #C and exponential w.r.t. M !
- Conditional independence assumptions :

$$p(\mathbf{c}|y) = \prod_{m=1}^{M} p(c_m|y).$$

- Linear memory complexity w.r.t. both #C and M.
- This is stacking, with Naive Bayes classifier as second-stage classifier !