## Machine Learning Techniques (機器學習技巧)



Lecture 7: Blending and Bagging

Hsuan-Tien Lin (林軒田) htlin@csie.ntu.edu.tw

Department of Computer Science & Information Engineering

National Taiwan University (國立台灣大學資訊工程系)



### Agenda

#### Lecture 7: Blending and Bagging

- Motivation of Aggregation
- Uniform Blending
- Linear and Any Blending
- Bagging

Motivation of Aggregation

### An Aggregation Story

Your *T* friends  $g_1, \dots, g_T$  predicts whether stock will go up as  $g_t(\mathbf{x})$ .

#### You can . . .

- select the most trust-worthy friend from their usual performance —validation!
- mix the predictions from all your friends uniformly —let them vote!
- mix the predictions from all your friends non-uniformly —let them vote, but give some more ballots
- combine the predictions conditionally
   —if [condition t true] give some ballots to friend t

# **aggregation** models: **mix** or **combine** hypotheses (for better performance)

Hsuan-Tien Lin (NTU CSIE)

. . .

Motivation of Aggregation

### Aggregation with Math Notations

Your *T* friends  $g_1, \dots, g_T$  predicts whether stock will go up as  $g_t(\mathbf{x})$ .

- select the most trust-worthy friend from their usual performance  $G(\mathbf{x}) = g_{t_*}(\mathbf{x})$  with  $t_* = \operatorname{argmin}_{t \in \{1, 2, \dots, T\}} E_{val}(g_t)$
- mix the predictions from all your friends uniformly  $Q(x) = x i m \left( \sum_{i=1}^{T} d_{i} x_{i} (x_{i}) \right)$

$$G(\mathbf{x}) = \operatorname{sign}\left(\sum_{t=1}^{T} \mathbf{1} \cdot g_t(\mathbf{x})\right)$$

- **mix** the predictions from all your friends **non-uniformly**  $G(\mathbf{x}) = \operatorname{sign}\left(\sum_{t=1}^{T} \alpha_t \cdot g_t(\mathbf{x})\right) \text{ with } \alpha_t \ge 0$ 
  - include select: α<sub>t</sub> = [[E<sub>val</sub>(g<sub>t</sub>) smallest]]
  - include uniformly:  $\alpha_t = 1$

### combine the predictions conditionally

 $G(\mathbf{x}) = \operatorname{sign}\left(\sum_{t=1}^{T} q_t(\mathbf{x}) \cdot g_t(\mathbf{x})\right) \text{ with } q_t(\mathbf{x}) \ge 0$ 

• include **non-uniformly**:  $q_t(\mathbf{x}) = \alpha_t$ 

#### aggregation models: a rich family

Hsuan-Tien Lin (NTU CSIE)

Motivation of Aggregation

Recall: Selection by Validation

$$G(\mathbf{x}) = g_{t_*}(\mathbf{x})$$
 with  $t_* = \operatorname*{argmin}_{t \in \{1, 2, \cdots, T\}} \frac{E_{\mathsf{val}}(g_t)}{E_{\mathsf{val}}(g_t)}$ 

- simple and popular
- can also use E<sub>in</sub> instead of E<sub>val</sub> (with complexity price on d<sub>VC</sub>)
- need one strong g<sub>t</sub> to guarantee small E<sub>val</sub> (and small E<sub>out</sub>)

selection: rely on one strong hypothesis aggregation: can we do better with many (possibly weaker) hypotheses?

Motivation of Aggregation

## Why Might Aggregation Work?



- mix different weak hypotheses uniformly —G(x) 'strong'



- mix different random-PLA hypotheses uniformly —G(x) 'moderate'
- aggregation → regularization (?)

proper aggregation  $\Longrightarrow$  better performance

Hsuan-Tien Lin (NTU CSIE)

Motivation of Aggregation

### Fun Time

# Uniform Blending (Voting) for Classification uniform blending: known $g_t$ , each with 1 ballot

$$G(\mathbf{x}) = \operatorname{sign}\left(\sum_{t=1}^{T} \mathbf{1} \cdot g_t(\mathbf{x})\right)$$

- same g<sub>t</sub> (autocracy): as good as one single g<sub>t</sub>
- very different g<sub>t</sub> (diversity + democracy): majority can correct minority
- similar results with uniform voting for multiclass

$$G(\mathbf{x}) = \operatorname*{argmax}_{1 \le k \le K} \sum_{t=1}^{T} \llbracket g_t(\mathbf{x}) = k \rrbracket$$



#### how about regression?

Hsuan-Tien Lin (NTU CSIE)

Uniform Blending

### Uniform Blending for Regression

 $G(\mathbf{x}) = \frac{1}{T} \sum_{t=1}^{T} g_t(\mathbf{x})$ 

- same g<sub>t</sub> (autocracy): as good as one single g<sub>t</sub>
- very different  $g_t$  (diversity + democracy):
  - some  $g_t(\mathbf{x}) > f(\mathbf{x})$ , some  $g_t(\mathbf{x}) < f(\mathbf{x})$
  - $\implies$  average **could be** more accurate than individual

## diverse hypotheses: even simple uniform blending can be better than one

### Theoretical Analysis of Uniform Blending

$$G(\mathbf{x}) = \frac{1}{T} \sum_{t=1}^{T} g_t(\mathbf{x})$$

$$avg((g_t(\mathbf{x}) - f(\mathbf{x}))^2) = avg(g_t^2 - 2g_t f + f^2)$$
  
=  $avg(g_t^2) - 2Gf + f^2$   
=  $avg(g_t^2) - G^2 + (G - f)^2$   
=  $avg(g_t^2) - 2G^2 + G^2 + (G - f)^2$   
=  $avg(g_t^2 - 2g_t G + G^2) + (G - f)^2$   
=  $avg((g_t - G)^2) + (G - f)^2$ 

$$\operatorname{avg}\left(E_{\operatorname{out}}(g_t)\right) = \operatorname{avg}\left(\mathcal{E}(g_t - G)^2\right) + E_{\operatorname{out}}(G)$$

Uniform Blending

### Some Special $g_t$

consider a virtual iterative process that for t = 1, 2, ..., T

- **1** request size-*N* data  $\mathcal{D}_t$  from  $P^N$  (i.i.d.)
- **2** obtain  $g_t$  by  $\mathcal{A}(\mathcal{D}_t)$

$$\bar{g} = \lim_{T \to \infty} G = \lim_{T \to \infty} \frac{1}{T} \sum_{t=1}^{T} g_t = \mathcal{E}_{\mathcal{D}} \mathcal{A}(\mathcal{D})$$

$$\operatorname{avg}\left(\mathcal{E}_{\operatorname{out}}(g_t)\right) = \operatorname{avg}\left(\mathcal{E}(g_t - \bar{g})^2\right) + \mathcal{E}_{\operatorname{out}}(\bar{g})$$

expected performance of A = expected deviation to consensus +performance of consensus

- performance of consensus: called bias
- expected deviation to consensus: called variance

#### uniform blending: reduces variance for stabler performance

Hsuan-Tien Lin (NTU CSIE)

### Fun Time

Linear and Any Blending

#### Linear Blending

linear blending: known  $g_t$ , each to be given  $\alpha_t$  ballot

$$G(\mathbf{x}) = \operatorname{sign}\left(\sum_{t=1}^{T} \alpha_t \cdot g_t(\mathbf{x})\right) \text{ with } \alpha_t \ge 0$$

computing 'good' 
$$\alpha_t$$
 :  $\min_{\alpha_t \ge 0} E_{in}(\alpha)$ 

 $\frac{\text{linear blending for regression}}{\min_{\alpha_t \ge 0} \frac{1}{N} \sum_{n=1}^{N} \left( y_n - \sum_{t=1}^{T} \alpha_t g_t(\mathbf{x}_n) \right)^2} \frac{\text{LinReg + transformation}}{\min_{w_i} \frac{1}{N} \sum_{n=1}^{N} \left( y_n - \sum_{i=1}^{\tilde{d}} w_i \phi_i(\mathbf{x}_n) \right)^2}$ 

linear blending = LinModel + hypotheses as transform + constraints

Hsuan-Tien Lin (NTU CSIE)

### Constraint on $\alpha_t$

linear blending = LinModel + hypotheses as transform + constraints:

$$\min_{t \ge 0} \qquad \frac{1}{N} \sum_{n=1}^{N} \operatorname{err} \left( y_n, \sum_{t=1}^{T} \alpha_t g_t(\mathbf{x}_n) \right)$$

linear blending for binary classification

$$\text{if } \alpha_t < 0 \implies \alpha_t g_t(\mathbf{x}) = |\alpha_t| \left( -g_t(\mathbf{x}) \right)$$

- negative  $\alpha_t$  for  $g_t \equiv$  positive  $|\alpha_t|$  for  $-g_t$
- if you have a stock up/down classifier with 99% error, tell me!
   :-)

in practice, often linear blending = LinModel + hypotheses as transform + constraints

Hsuan-Tien Lin (NTU CSIE)

### Linear Blending versus Selection

in practice, often

$$\textbf{g}_1 \in \mathcal{H}_1, \textbf{g}_2 \in \mathcal{H}_2, \dots, \textbf{g}_T \in \mathcal{H}_T$$

by minimum E<sub>in</sub>

- recall: selection by minimum  $E_{in}$ —best of best, paying  $d_{VC} \left( \bigcup_{t=1}^{T} \mathcal{H}_{t} \right)$
- recall: linear blending includes selection as special case
   —by setting α<sub>t</sub> = [[E<sub>val</sub>(g<sub>t</sub>) smallest]]
- complexity price of linear blending with  $E_{in}$  (aggregation of best):  $\gg d_{VC} \left( \bigcup_{t=1}^{T} \mathcal{H}_t \right)$

# like selection, blending practically done with $(E_{val} \text{ instead of } E_{in}) + (g_t \text{ from } E_{train})$

Hsuan-Tien Lin (NTU CSIE)

## Any Blending

#### Linear Blending

Given  $g_1, g_2, ..., g_T$ 

1 transform  $(\mathbf{x}_n, y_n)$  in  $\mathcal{D}_{val}$  to  $(\mathbf{z}_n = \mathbf{\Phi}(\mathbf{x}_n), y_n)$ , where  $\mathbf{\Phi}(\mathbf{x}) = (g_1(\mathbf{x}), \dots, g_T(\mathbf{x}))$ 

**2** compute 
$$\alpha = \text{Lin}(\{(\mathbf{z}_n, y_n)\})$$

return  $G_{\text{LINB}}(\mathbf{x}) = \text{LinH}(\boldsymbol{\alpha}^T \boldsymbol{\Phi}(\mathbf{x}))$ 

#### Any Blending (Stacking)

Given  $g_1, g_2, ..., g_T$ 

1 transform  $(\mathbf{x}_n, y_n)$  in  $\mathcal{D}_{val}$  to  $(\mathbf{z}_n = \mathbf{\Phi}(\mathbf{x}_n), y_n)$ , where  $\mathbf{\Phi}(\mathbf{x}) = (g_1(\mathbf{x}), \dots, g_T(\mathbf{x}))$ 2 compute  $\tilde{g} = \operatorname{Any}(\{(\mathbf{z}_n, y_n)\})$ 

return  $G_{\text{ANYB}}(\mathbf{x}) = \tilde{g}(\mathbf{\Phi}(\mathbf{x}))$ 

if AnyModel = quadratic polynomial:

$$G_{\text{ANYB}}(\mathbf{x}) = \sum_{t=1}^{T} \underbrace{\left(\alpha_t + \sum_{\tau=1}^{T} \alpha_{\tau,t} g_{\tau}(\mathbf{x})\right)}_{q(\mathbf{x})} \cdot g_t(\mathbf{x}) - \text{conditional aggregation}$$

#### danger: overfitting with any blending!

Hsuan-Tien Lin (NTU CSIE)



KDDCup 2012 Track 1: World Champion Solution by NTU

• validation set blending: a special any blending model

 $E_{\text{test}}$  (squared): 519.45  $\implies$  456.24

-helped secure the lead in last two weeks

test set blending: linear blending using *E*<sub>test</sub>

 $E_{\text{test}}$  (squared): 456.24  $\Longrightarrow$  442.06

-helped turn the tables in last hour

## blending 'useful' in practice, despite the computational burden

Hsuan-Tien Lin (NTU CSIE)

Linear and Any Blending

### Fun Time

Hsuan-Tien Lin (NTU CSIE)

Blending and Bagging	Bagging		
	What We Have Done		
	blending: aggregate after getting $g_t$ ;		
	icarning. aggregate as well as getting gr		
	aggregation type	blending	learning
	uniform	voting/averaging	?
	non-uniform	linear	?
	conditional	stacking	?

learning  $g_t$  for uniform aggregation: diversity important

- diversity by different models:  $g_1 \in \mathcal{H}_1, g_2 \in \mathcal{H}_2, \dots, g_T \in \mathcal{H}_T$
- diversity by different parameters: GD with  $\eta = 0.001, 0.01, ..., 10$
- diversity by algorithmic randomness: random PLA with different random seeds
- diversity by data randomness:

within-cross-validation hypotheses  $g_v^-$ 

next: diversity by data randomness without  $g^-$ 

Hsuan-Tien Lin (NTU CSIE)

### **Revisit of Bias-Variance**

expected performance of A = expected deviation to consensus +performance of consensus

consensus  $\bar{g} = \text{expected } g_t \text{ from } \mathcal{D}_t \sim \mathcal{P}^N$ 

- consensus more stable than direct A(D), but comes from many more D<sub>t</sub> than the D on hand
- want: approximate  $\bar{g}$  by
  - finite (large) T
  - approximate  $g_t = \mathcal{A}(\mathcal{D}_t)$  from  $\mathcal{D}_t \sim P^N$  using only  $\mathcal{D}$

## bootstrapping: a statistical tool that re-samples from $\mathcal{D}$ to 'simulate' $\mathcal{D}_t$

## Bootstrap Aggregation

#### bootstrapping

bootstrap sample  $\tilde{D}_t$ : re-sample N examples from  $\mathcal{D}$  with replacement

#### virtual aggregation

consider a **virtual** iterative process that for t = 1, 2, ..., T

1 request size-*N* data  $D_t$  from  $P^N$  (i.i.d.)

2 obtain 
$$g_t$$
 by  $\mathcal{A}(\mathcal{D}_t)$ 

 $G = Uniform(g_t)$ 

#### bootstrap aggregation

consider a **physical** iterative process that for t = 1, 2, ..., T

1 request size-*N* data  $\tilde{\mathcal{D}}_t$  from bootstrapping

**2** obtain 
$$g_t$$
 by  $\mathcal{A}(\tilde{\mathcal{D}}_t)$ 

 $G = \text{Uniform}(g_t)$ 

bootstrap aggregation (BAGging): a simple meta algorithm on top of base algorithm  $\mathcal{A}$ 

Hsuan-Tien Lin (NTU CSIE)

#### Bagging

## **Bagging Pocket in Action**



 $T_{\text{pocket}} = 1000; \ T_{\text{bag}} = 25$ 

- very diverse *g*<sub>t</sub> from bagging
- proper non-linear boundary after aggregating binary classifiers

#### bagging works reasonably well if base algorithm sensitive to data randomness

Hsuan-Tien Lin (NTU CSIE)

Bagging

### Fun Time

Hsuan-Tien Lin (NTU CSIE)

Machine Learning Techniques

22/23

### Summary

#### Lecture 7: Blending and Bagging

Motivation of Aggregation

#### strong and/or moderate

- Uniform Blending one hypothesis, one vote, one value
- Linear and Any Blending learning with hypotheses as transform
- Bagging

#### bootstrapping for diverse $g_t$