

Ensemble

團隊合作，一把model一起上！

Framework of Ensemble

有Bagging & Boosting 場合不一樣

- Get a set of classifiers

- $f_1(x), f_2(x), f_3(x), \dots$



坦



補



DD

They should be diverse.

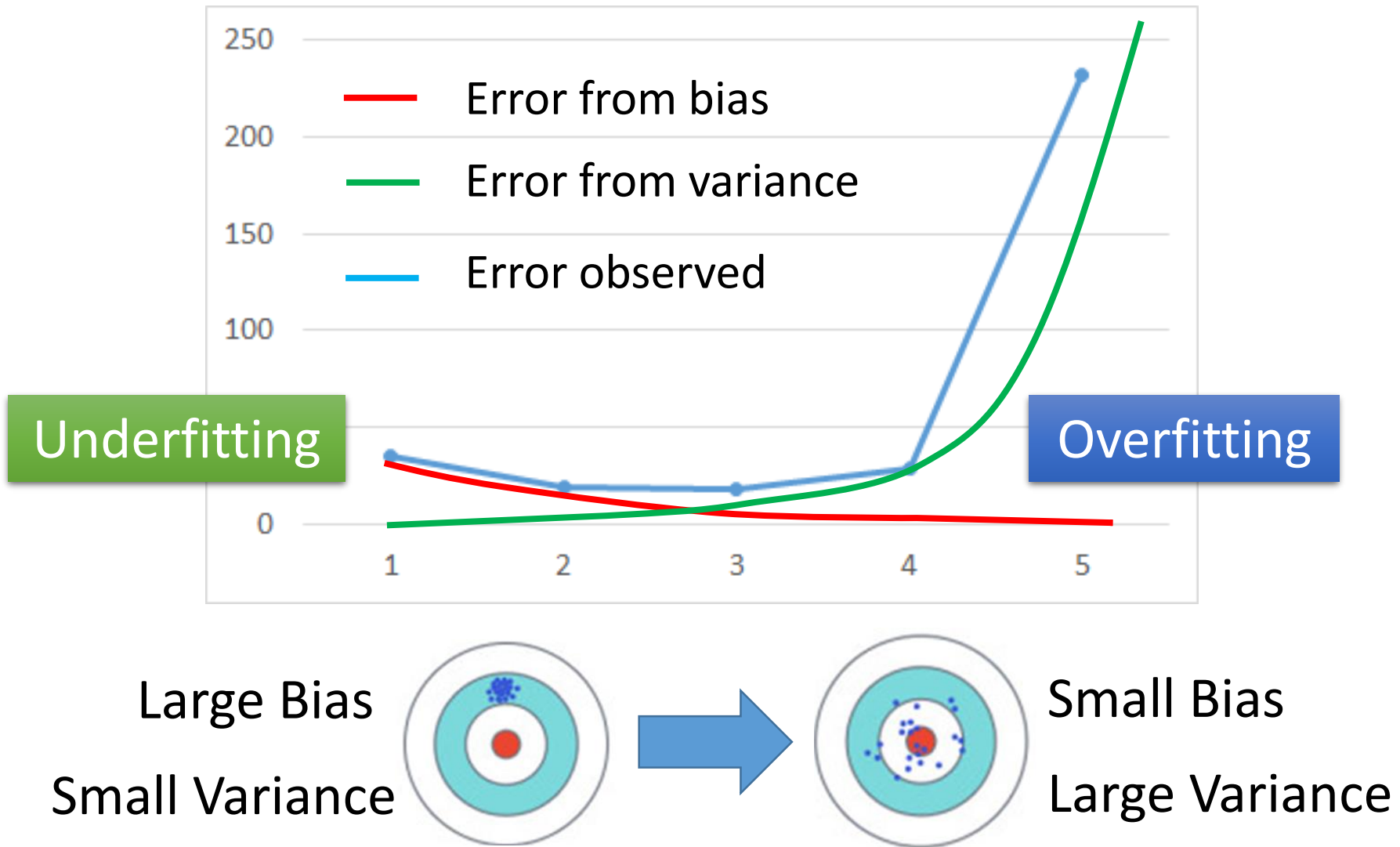
整合

- Aggregate the classifiers (*properly*)
 - 在打王時每個人都有該站的位置

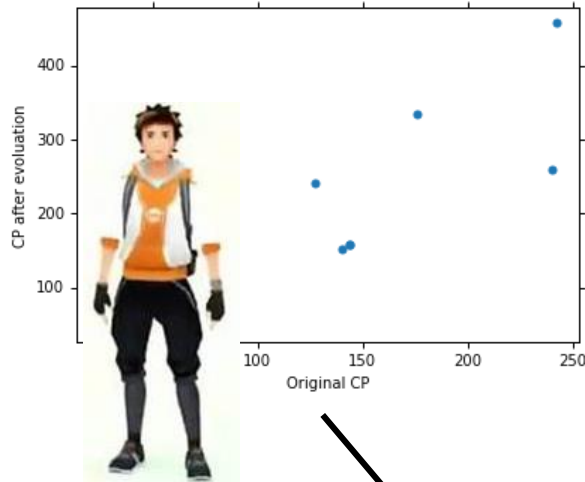
Ensemble: Bagging

Bagging中model的訓練是沒有順序的，可以一次train一坨

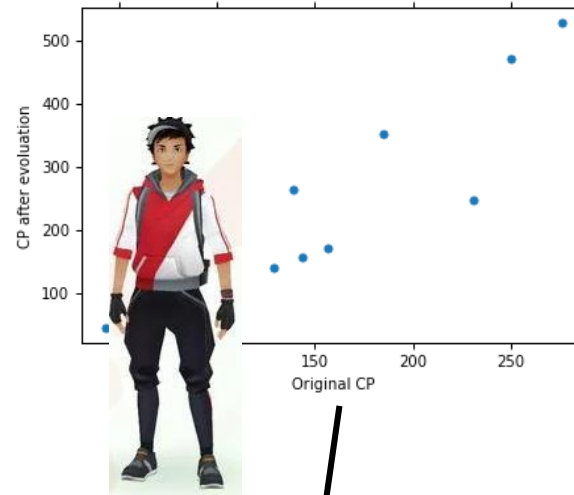
Review: Bias v.s. Variance



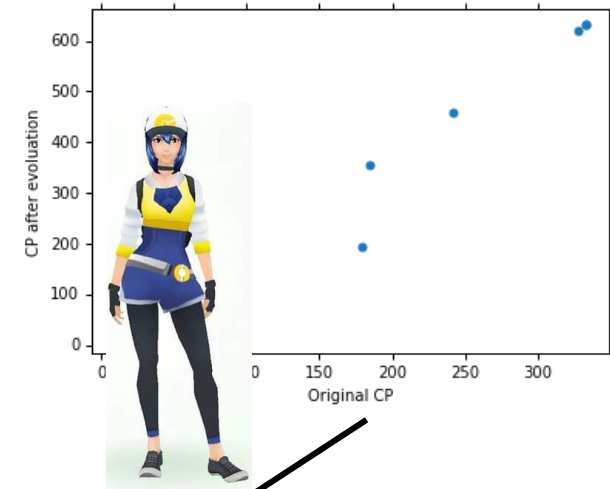
Universe 1



Universe 2

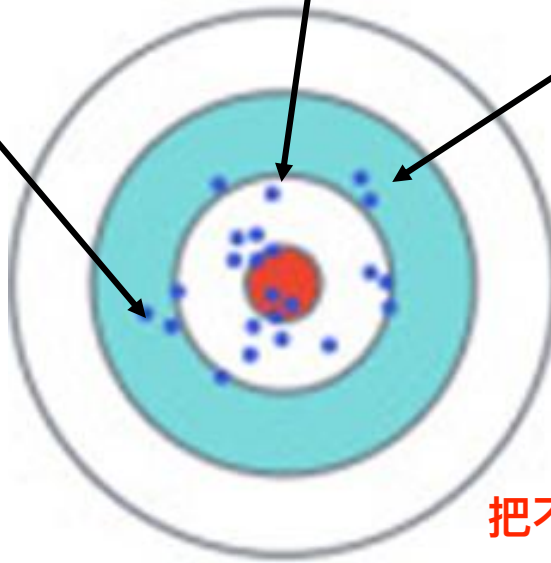


Universe 3



A complex model will have large variance.

We can average complex models to reduce variance.

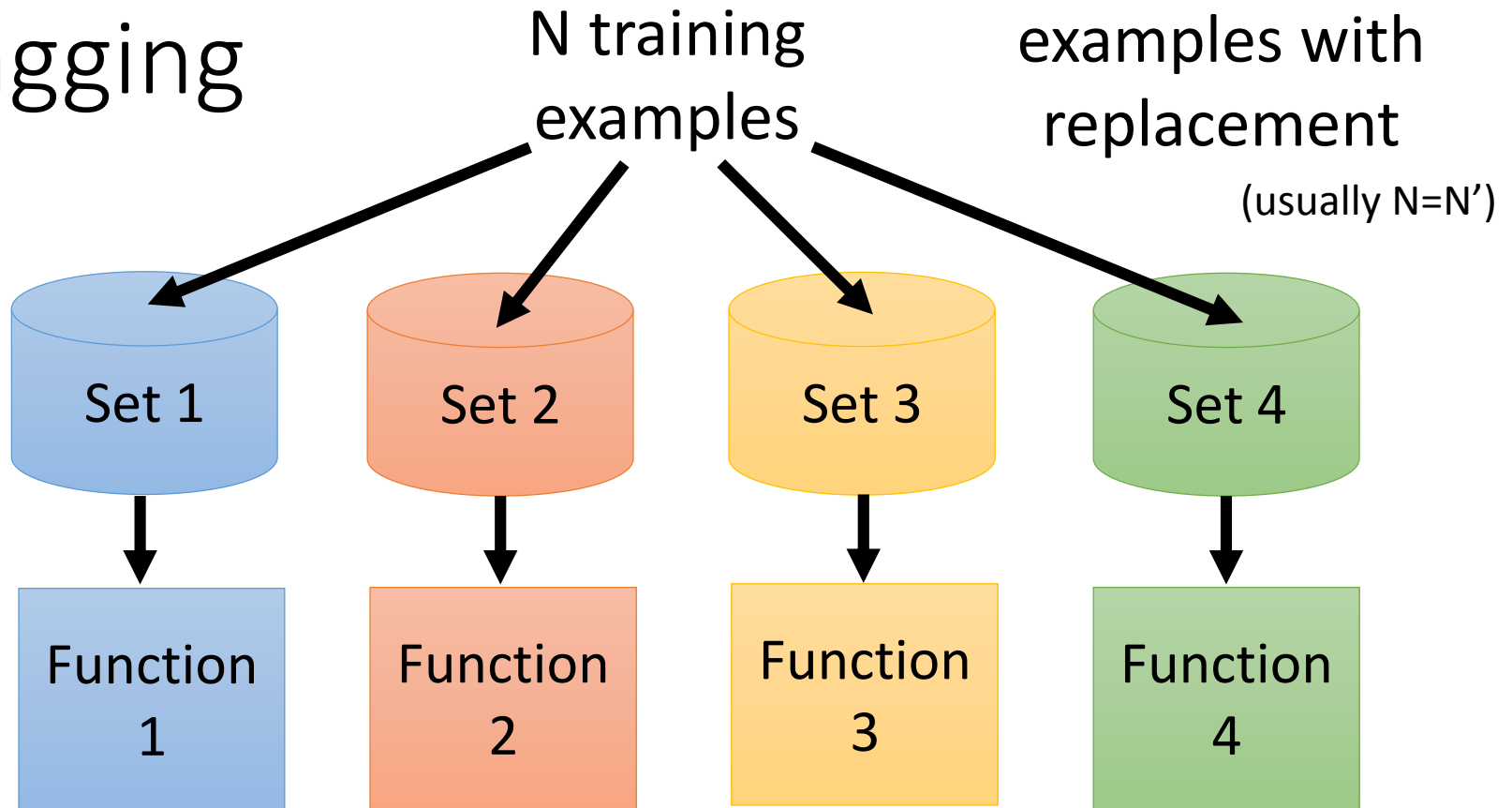


If we average all the f^* , is it close to \hat{f}

$$E[f^*] = \hat{f}$$

把不同的模型的輸出集合起來做平均

Bagging



把training data先分成很多份，在各自針對sub training data去train出一個model

Bagging

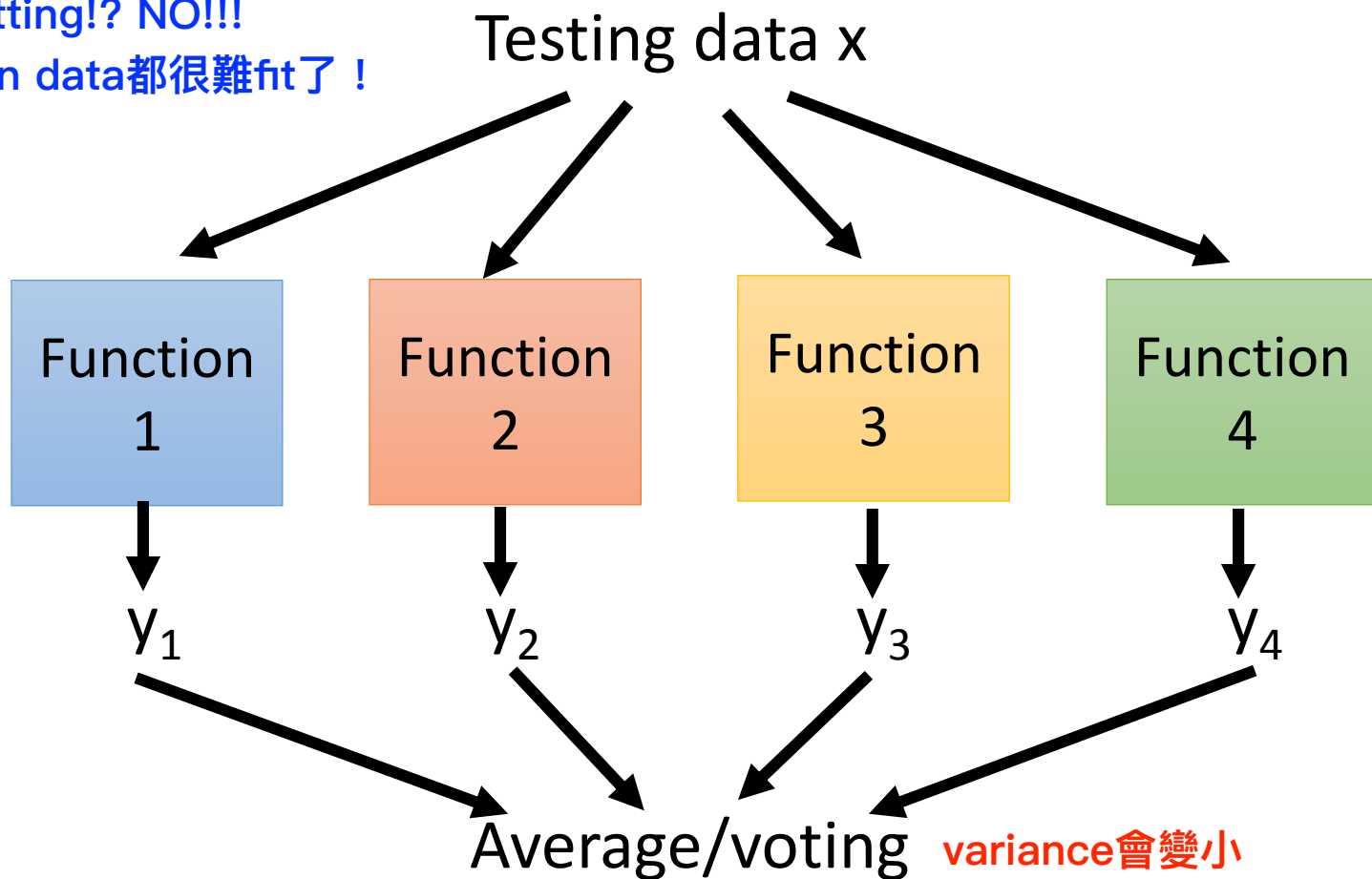
This approach would be helpful when your model is complex, easy to overfit.

當model很複雜，擔心他overfitting的時候，為了減低variance才做bagging

e.g. decision tree

NN overfitting!? NO!!!

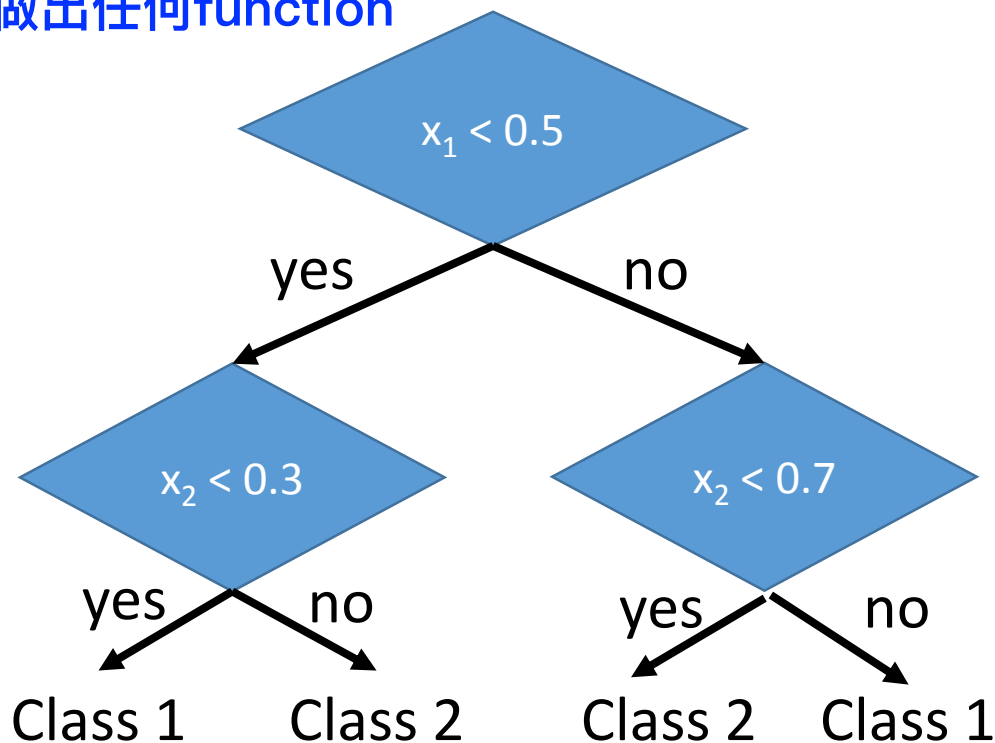
通常在train data都很難fit了！



decision tree非常容易overfitting，
因此可做bagging，結果即為
random forest

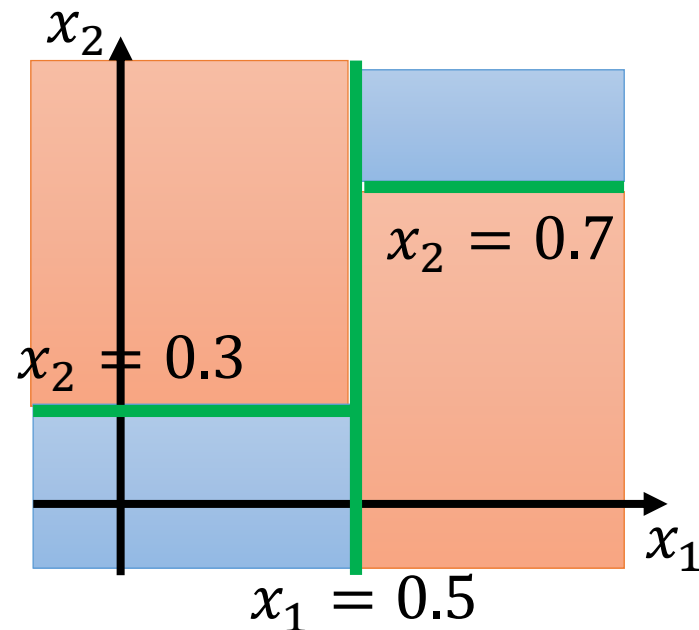
Decision Tree

樹夠深，decision tree
可做出任何function



Can have more complex questions

Assume each object x is
represented by a 2-dim vector $\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$



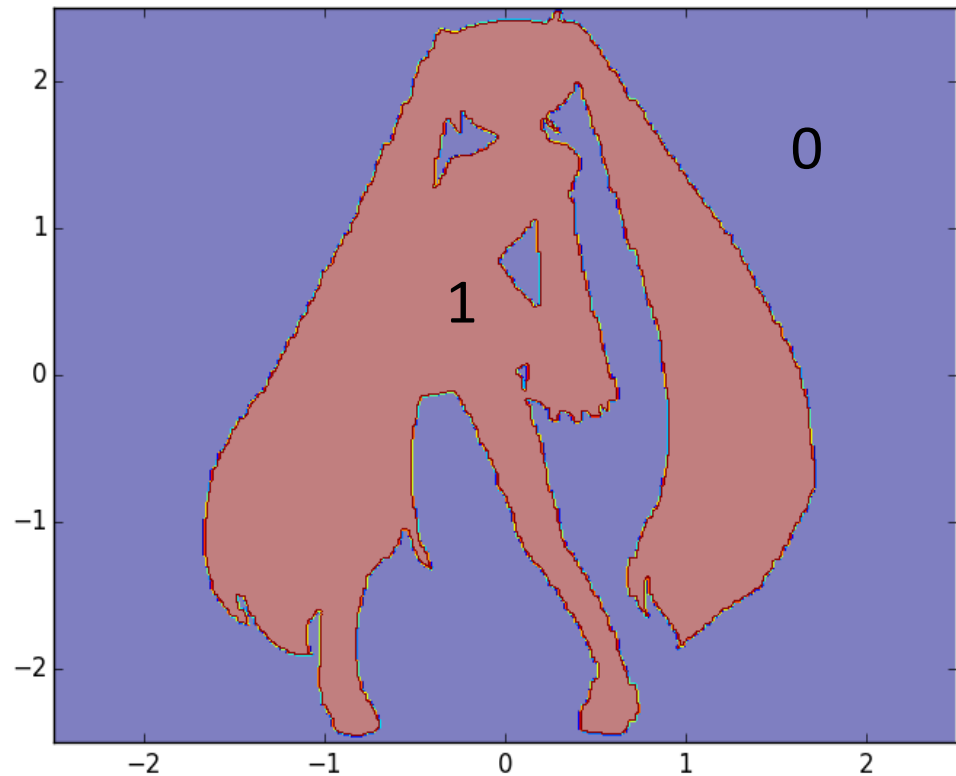
The questions in
training

number of branches,
Branching criteria,
termination criteria,

每個節點要做多少分支，criteria，什麼時候停止分支..等等需要考慮 base hypothesis

decision tree example

Experiment: Function of Miku

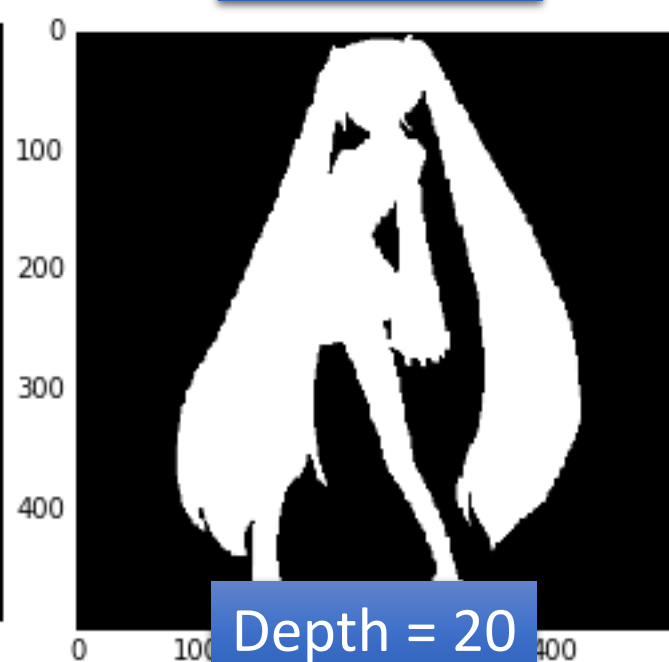
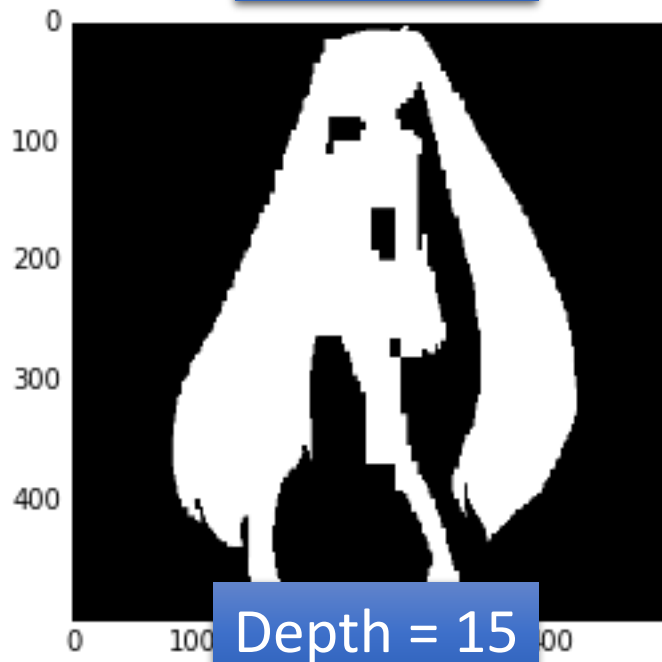
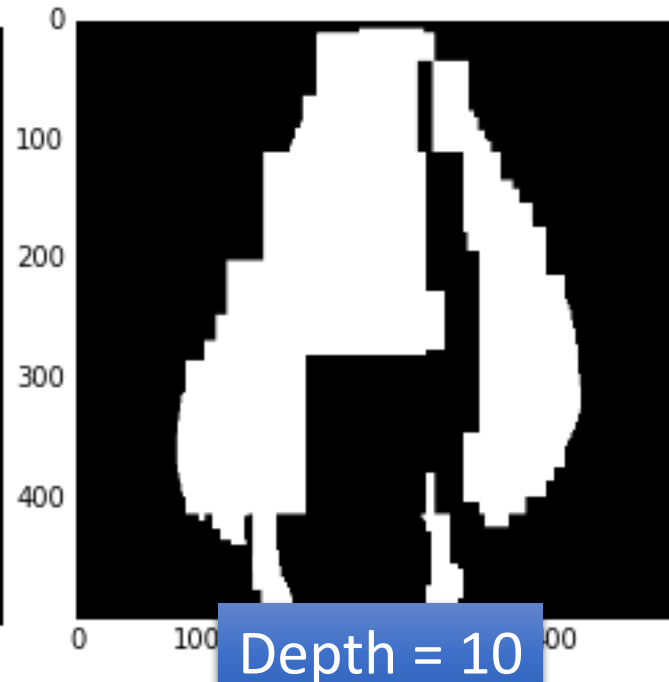
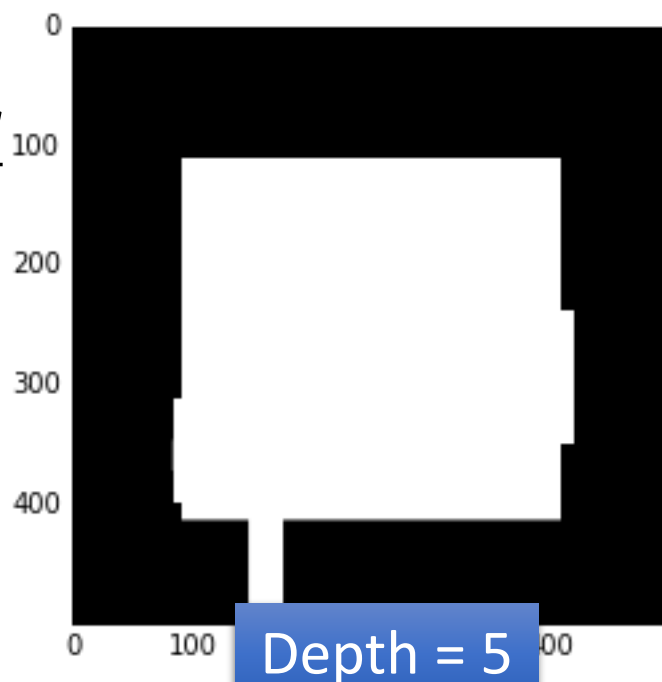


http://speech.ee.ntu.edu.tw/~tlkagk/courses/MLDS_2015_2/theano/miku

(1st column: x, 2nd column: y, 3rd column: output (1 or 0))

Experiment:
Function of Miku

Single
Decision
Tree



Random Forest

train	f_1	f_2	f_3	f_4
x^1	O	X	O	X
x^2	O	X	X	O
x^3	X	O	O	X
x^4	X	O	X	O

- Decision tree:
 - Easy to achieve 0% error rate on training data
 - If each training example has its own leaf
 - Random forest: Bagging of decision tree
 - Resampling training data is not sufficient
 - Randomly restrict the features/questions used in each split random決定哪些feature不能用 特別的sample方法
 - Out-of-bag validation for bagging 不需要切validation也可以有validation的效果
 - Using RF = $f_2 + f_4$ to test x^1
 - Using RF = $f_2 + f_3$ to test x^2
 - Using RF = $f_1 + f_4$ to test x^3
 - Using RF = $f_1 + f_3$ to test x^4

Out-of-bag (OOB) error 取平均

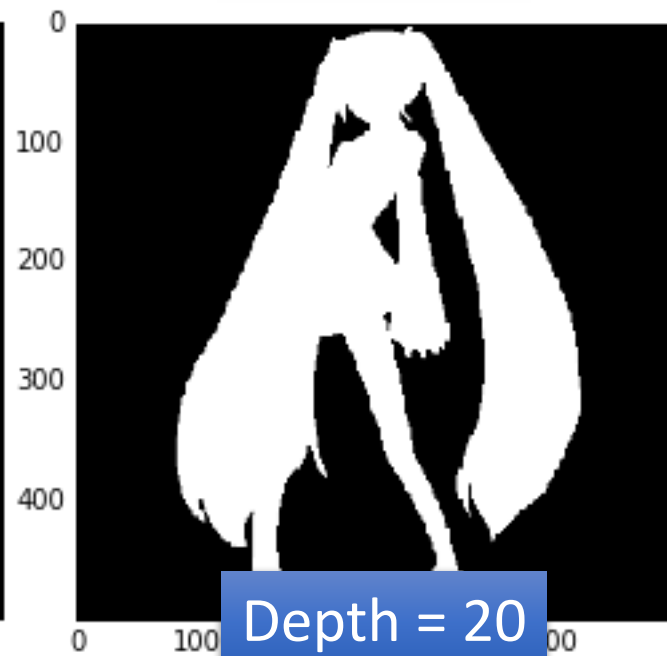
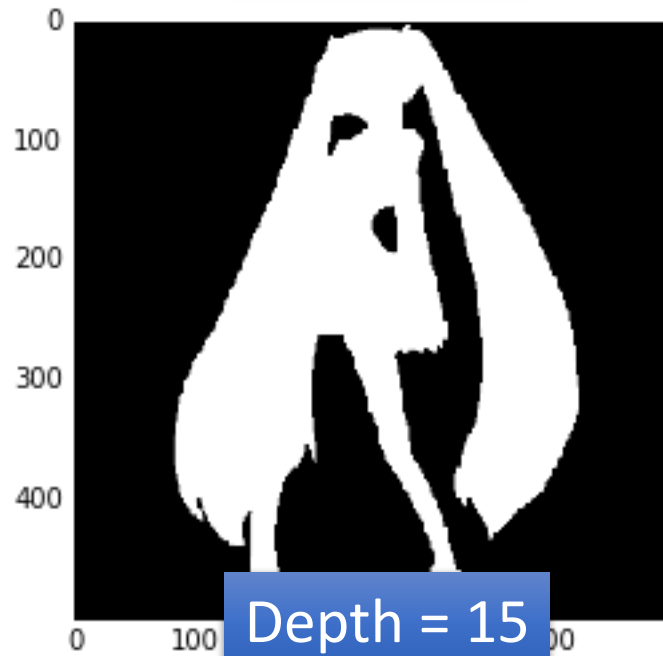
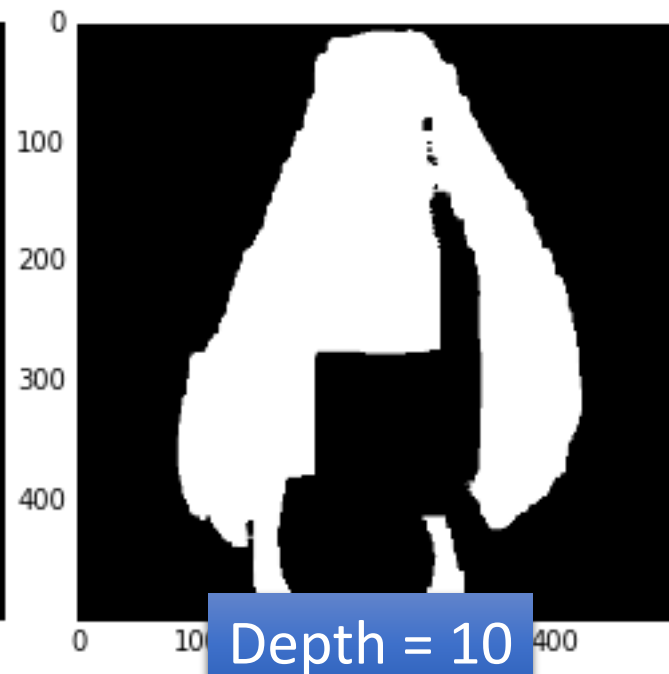
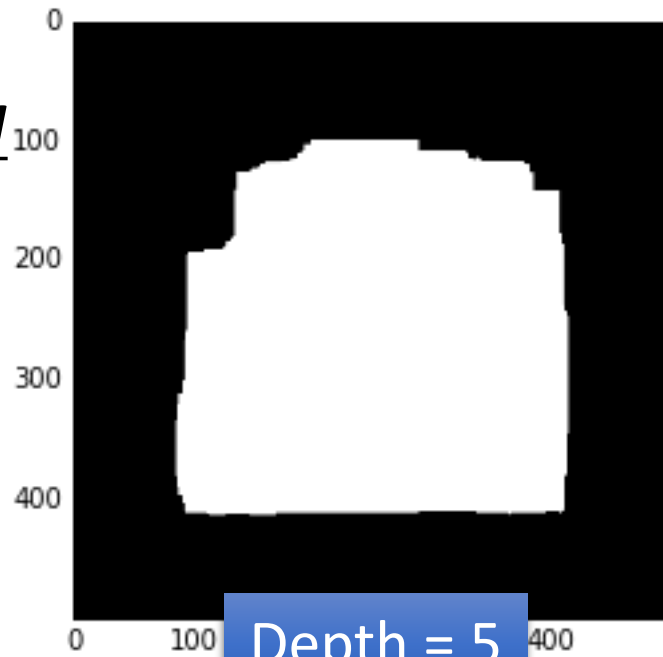
Good error estimation of testing set
- 拿沒看過test data的model去test他的validation accuracy

Experiment: Function of Miku

用random forest並不會讓你更容易去fit data，只是讓結果比較平滑 (robust)

Random
Forest

(100 trees)



Ensemble: Boosting

Improving Weak Classifiers

bagging是用在很強的model，Boosting是用在很弱的model，當有些raw model但是無法讓他們fit data，這時候可以採用boosting

Boosting

Training data:

$$\{(x^1, \hat{y}^1), \dots, (x^n, \hat{y}^n), \dots, (x^N, \hat{y}^N)\}$$

$\hat{y} = \pm 1$ (binary classification)

- Guarantee: 神保證！！
 - If your ML algorithm can produce classifier with error rate smaller than 50% on training data
 - You can obtain 0% error rate classifier after boosting.
- Framework of boosting
 - Obtain the first classifier $f_1(x)$
 - Find another function $f_2(x)$ to help $f_1(x)$
 - However, if $f_2(x)$ is similar to $f_1(x)$, it will not help a lot. f2雖然是輔助f1，但是他們彼此需要是互補的
 - We want $f_2(x)$ to be complementary with $f_1(x)$ (How?)
 - Obtain the second classifier $f_2(x)$
 - Finally, combining all the classifiers
- The classifiers are learned sequentially.

model的訓練是有順序的，比須先找出f1，在找出與其互補的f2...以此類推

How to obtain different classifiers?

- Training on different training data sets
- How to have different training data sets
 - Re-sampling your training data to form a new set
 - Re-weighting your training data to form a new set
 - In real implementation, you only have to change the cost/objective function

$$(x^1, \hat{y}^1, u^1) \quad u^1 = \cancel{1} \quad 0.4$$

$$(x^2, \hat{y}^2, u^2) \quad u^2 = \cancel{1} \quad 2.1$$

$$(x^3, \hat{y}^3, u^3) \quad u^3 = \cancel{1} \quad 0.7$$

$$L(f) = \sum_n l(f(x^n), \hat{y}^n)$$



$$L(f) = \sum_n \boxed{u^n} l(f(x^n), \hat{y}^n)$$

乘上權重

藉由改變weight來產生不同的data set

sampling也可以想像成改weight，只是weight變成整數（sample到的次數）

boosting的其中一種方法

Idea of Adaboost

training example # : n

step 2. 找一組新的training data在f1上面是會壞掉的

- Idea: training $f_2(x)$ on the new training set that fails $f_1(x)$
- How to find a new training set that fails $f_1(x)$?

ε_1 : the error rate of $f_1(x)$ on its training data

$$\varepsilon_1 = \frac{\sum_n \overset{\text{weight}}{u_1^n} \delta(f_1(x^n) \neq \hat{y}^n)}{Z_1 \text{ normalization}} \quad Z_1 = \sum_n u_1^n \quad \begin{matrix} \text{假設錯誤率小於五十趴} \\ \varepsilon_1 < 0.5 \end{matrix}$$

step 2.

Changing the example weights from u_1^n to u_2^n such that

$$\frac{\sum_n u_2^n \delta(f_1(x^n) \neq \hat{y}^n)}{Z_2} = 0.5$$

使其壞掉

The performance of f_1 for new weights would be random.

Training $f_2(x)$ based on the new weights u_2^n

Re-weighting Training Data

- Idea: training $f_2(x)$ on the new training set that fails $f_1(x)$
- How to find a new training set that fails $f_1(x)$?

把答對的題目權重變小，把答錯的題目權重變大 X D 這樣就可以把 f_1 的結果搞爛找出一組與其互補的 new data set

$$(x^1, \hat{y}^1, u^1) \quad u^1 = 1 \quad \checkmark$$

$$(x^2, \hat{y}^2, u^2) \quad u^2 = 1 \quad \times$$

$$(x^3, \hat{y}^3, u^3) \quad u^3 = 1 \quad \checkmark$$

$$(x^4, \hat{y}^4, u^4) \quad u^4 = 1 \quad \checkmark$$

$$u^1 = 1/\sqrt{3} \quad \checkmark$$

$$u^2 = \sqrt{3} \quad \times$$

$$u^3 = 1/\sqrt{3} \quad \checkmark$$

$$u^4 = 1/\sqrt{3} \quad \checkmark$$

$$\varepsilon_1 = 0.25$$

$$1/4$$

$f_1(x)$

$$0.5$$

$$\varepsilon_2 < 0.5$$

$f_2(x)$

拿新的 data set 去 train 出 f_2 ，訓練出 error rate < 0.5

Re-weighting Training Data

- Idea: **training $f_2(x)$ on the new training set that fails $f_1(x)$**
- How to find a new training set that fails $f_1(x)$?

$$\left\{ \begin{array}{ll} \text{If } x^n \text{ misclassified by } f_1 \text{ (} f_1(x^n) \neq \hat{y}^n \text{)} & u_2^n \leftarrow u_1^n \text{ multiplying } d_1 \quad \text{increase} \\ \text{If } x^n \text{ correctly classified by } f_1 \text{ (} f_1(x^n) = \hat{y}^n \text{)} & u_2^n \leftarrow u_1^n \text{ divided by } d_1 \quad \text{decrease} \end{array} \right.$$

f_2 will be learned based on example weights u_2^n

What is the value of d_1 ?

Re-weighting Training Data

$$\varepsilon_1 = \frac{\sum_n u_1^n \delta(f_1(x^n) \neq \hat{y}^n)}{Z_1}$$

$$Z_1 = \sum_n u_1^n$$

分類錯誤的地方

$$\frac{\sum_n u_2^n \delta(f_1(x^n) \neq \hat{y}^n)}{Z_2} = 0.5$$

$$f_1(x^n) \neq \hat{y}^n \quad u_2^n \leftarrow u_1^n \text{ multiplying } d_1$$

$$f_1(x^n) = \hat{y}^n \quad u_2^n \leftarrow u_1^n \text{ divided by } d_1$$

$$= \sum_{f_1(x^n) \neq \hat{y}^n} u_1^n d_1$$

$$= \sum_{f_1(x^n) \neq \hat{y}^n} u_2^n + \sum_{f_1(x^n) = \hat{y}^n} u_2^n$$

$$= \sum_n u_2^n$$

$$= \sum_{f_1(x^n) \neq \hat{y}^n} u_1^n d_1 + \sum_{f_1(x^n) = \hat{y}^n} u_1^n / d_1$$

分類錯誤的地方

分類正確的地方

$$\frac{\sum_{f_1(x^n) \neq \hat{y}^n} u_1^n d_1 + \sum_{f_1(x^n) = \hat{y}^n} u_1^n / d_1}{\sum_{f_1(x^n) \neq \hat{y}^n} u_1^n d_1} = 2$$

Re-weighting Training Data

epsilon <= 0.5

$$\varepsilon_1 = \frac{\sum_n u_1^n \delta(f_1(x^n) \neq \hat{y}^n)}{Z_1}$$

$$Z_1 = \sum_n u_1^n$$

$$\frac{\sum_n u_2^n \delta(f_1(x^n) \neq \hat{y}^n)}{Z_2} = 0.5$$

$$\begin{aligned} f_1(x^n) \neq \hat{y}^n & \quad u_2^n \leftarrow u_1^n \text{ multiplying } d_1 \\ f_1(x^n) = \hat{y}^n & \quad u_2^n \leftarrow u_1^n \text{ divided by } d_1 \end{aligned}$$

$$\frac{\sum_{f_1(x^n) \neq \hat{y}^n} u_1^n d_1 + \sum_{f_1(x^n) = \hat{y}^n} u_1^n / d_1}{\sum_{f_1(x^n) \neq \hat{y}^n} u_1^n d_1} = 2 \quad \frac{\sum_{f_1(x^n) = \hat{y}^n} u_1^n / d_1}{\sum_{f_1(x^n) \neq \hat{y}^n} u_1^n d_1} = 1$$

$$\sum_{f_1(x^n) = \hat{y}^n} u_1^n / d_1 = \sum_{f_1(x^n) \neq \hat{y}^n} u_1^n d_1$$

$$\varepsilon_1 = \frac{\sum_{f_1(x^n) \neq \hat{y}^n} u_1^n}{Z_1} \quad \text{答錯的總和}$$

$$\sum_{f_1(x^n) \neq \hat{y}^n} u_1^n = Z_1 \varepsilon_1$$

$$\frac{1}{d_1} \sum_{f_1(x^n) = \hat{y}^n} u_1^n = d_1 \sum_{f_1(x^n) \neq \hat{y}^n} u_1^n$$

$$\frac{Z_1(1 - \varepsilon_1)}{d_1} = Z_1 \varepsilon_1 d_1$$

$$Z_1(1 - \varepsilon_1)/d_1 = Z_1 \varepsilon_1 d_1$$

$$d_1 = \sqrt{(1 - \varepsilon_1)/\varepsilon_1} > 1 \quad \text{結論}$$

Algorithm for AdaBoost

- Giving training data $\{(x^1, \hat{y}^1, u_1^1), \dots, (x^n, \hat{y}^n, u_1^n), \dots, (x^N, \hat{y}^N, u_1^N)\}$
 - $\hat{y} = \pm 1$ (Binary classification), $u_1^n = 1$ (equal weights)
- For $t = 1, \dots, T$:
 - Training weak classifier $f_t(x)$ with weights $\{u_t^1, \dots, u_t^N\}$
 - ε_t is the error rate of $f_t(x)$ with weights $\{u_t^1, \dots, u_t^N\}$
 - For $n = 1, \dots, N$:

分類錯誤

- If x^n is misclassified by $f_t(x)$: $\hat{y}^n \neq f_t(x^n)$

- $u_{t+1}^n = u_t^n \times d_t = u_t^n \times \exp(\alpha_t) \quad d_t = \sqrt{(1 - \varepsilon_t) / \varepsilon_t}$

分類正確

- Else:
- $u_{t+1}^n = u_t^n / d_t = u_t^n \times \exp(-\alpha_t) \quad \alpha_t = \ln \sqrt{(1 - \varepsilon_t) / \varepsilon_t}$

$$u_{t+1}^n \leftarrow u_t^n \times \exp(-\hat{y}^n * f_t(x^n) * \alpha_t)$$

換成exponential是為了將式子表達簡略

Algorithm for AdaBoost

- We obtain a set of functions: $f_1(x), \dots, f_t(x), \dots, f_T(x)$
- How to aggregate them?
 - Uniform weight:
 - $H(x) = \text{sign}(\sum_{t=1}^T f_t(x))$
 - Non-uniform weight:
 - $H(x) = \text{sign}(\sum_{t=1}^T \alpha_t f_t(x))$

Smaller error ε_t ,
larger weight for
final voting

因為每個classification有不同的好壞之分，因此需要在前面先呈上一個權重 (alpha t)

$$\alpha_t = \ln \sqrt{(1 - \varepsilon_t) / \varepsilon_t}$$

$$\varepsilon_t = 0.1$$

$$\varepsilon_t = 0.4$$

$$u_{t+1}^n = u_t^n \times \exp(-\hat{y}^n f_t(x^n) \alpha_t)$$

$$\alpha_t = 1.10$$

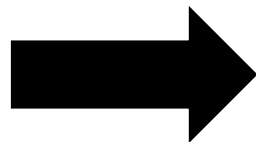
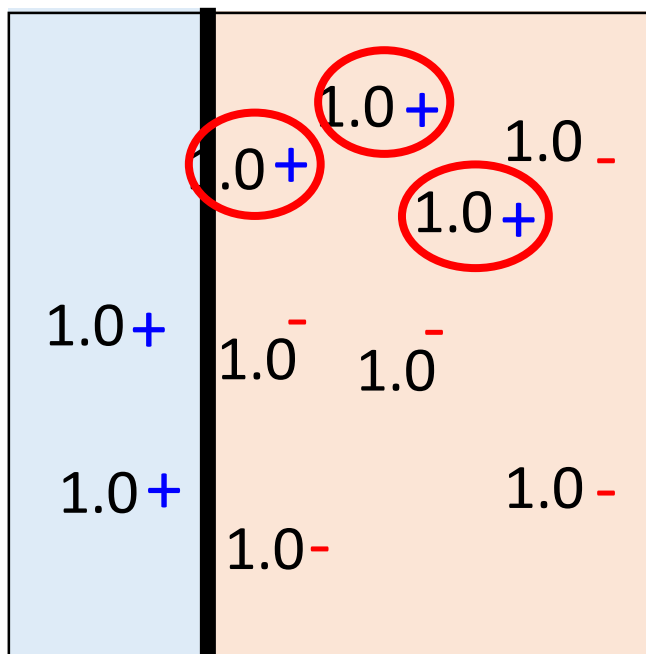
$$\alpha_t = 0.20$$

Toy Example

假設分佈在二維平面上，切一刀分類（超弱）

$T=3$, weak classifier = decision stump

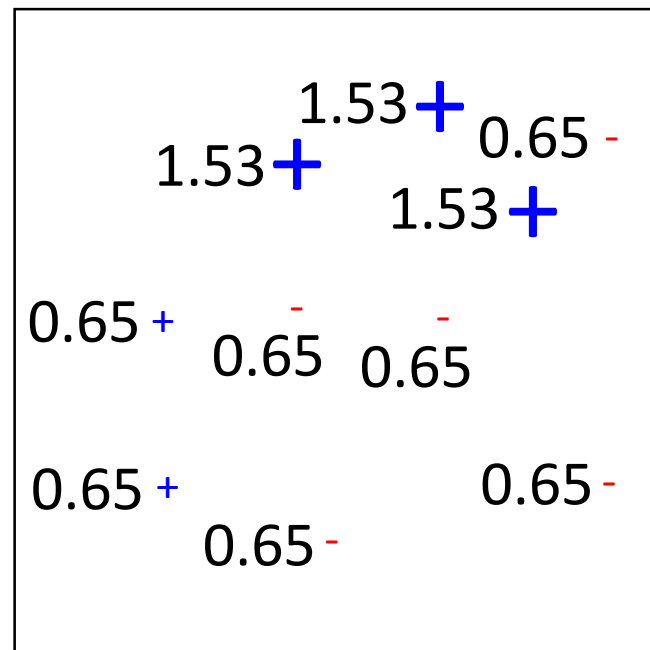
- $t=1$



$$\varepsilon_1 = 0.30$$

$$d_1 = 1.53$$

$$\alpha_1 = 0.42$$



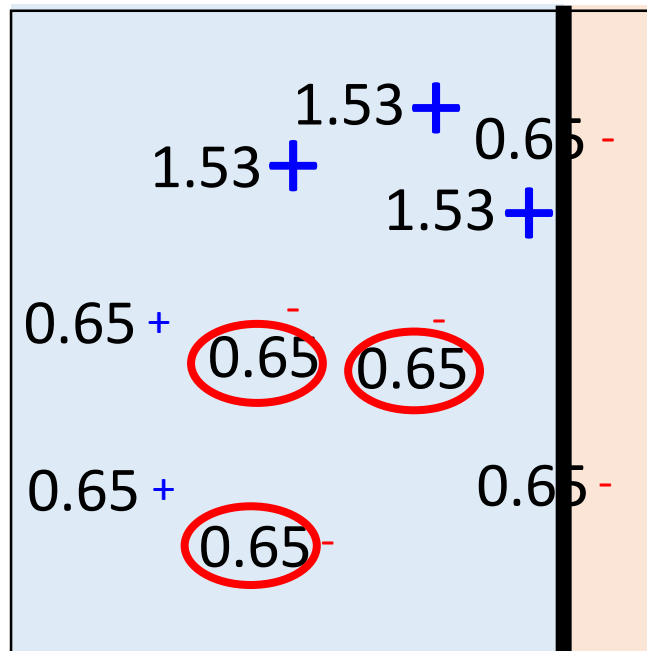
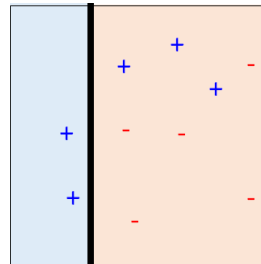
Toy Example

$T=3$, weak classifier = decision stump

• $t=2$

$$f_1(x):$$

$$\alpha_1 = 0.42$$



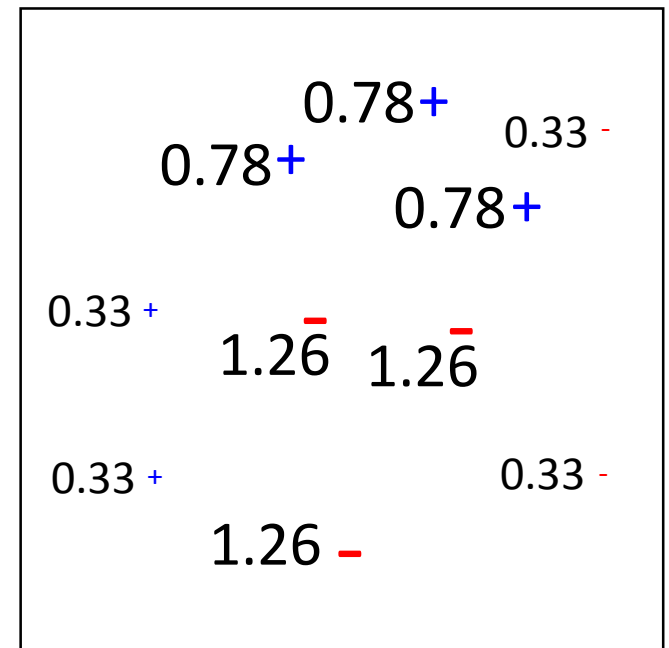
$$f_2(x)$$



$$\varepsilon_2 = 0.21$$

$$d_2 = 1.94$$

$$\alpha_2 = 0.66$$



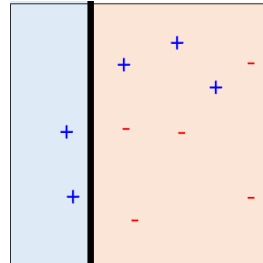
Toy Example

$T=3$, weak classifier = decision stump

• $t=3$

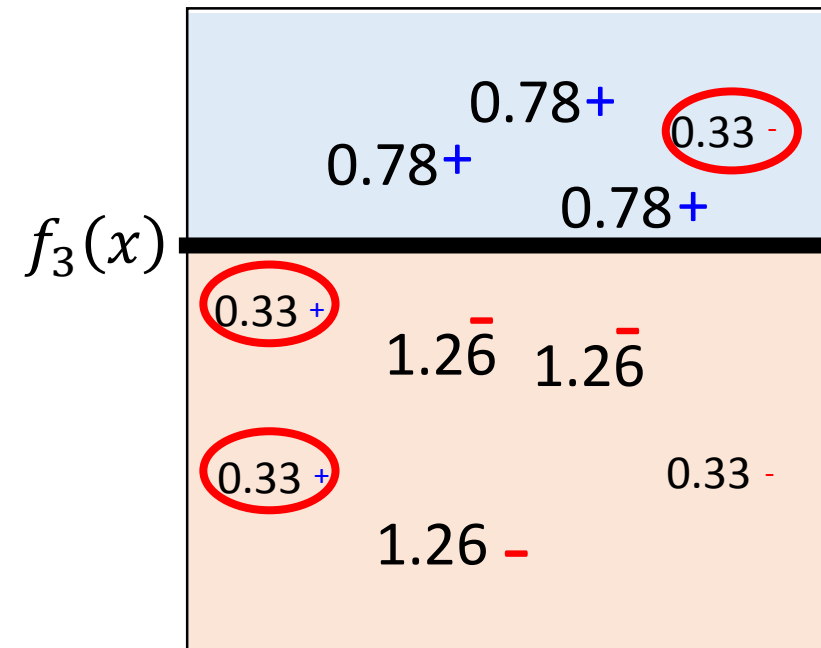
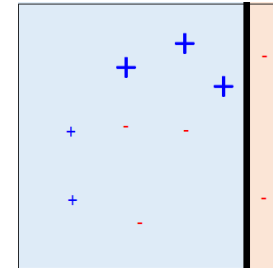
$f_1(x)$:

$$\alpha_1 = 0.42$$



$f_2(x)$:

$$\alpha_2 = 0.66$$



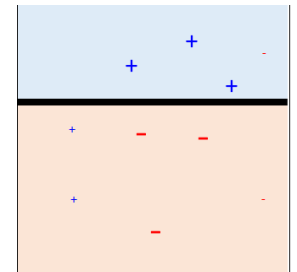
$$\varepsilon_3 = 0.13$$

$$d_3 = 2.59$$

$$\alpha_3 = 0.95$$

$f_3(x)$:

$$\alpha_3 = 0.95$$



Toy Example

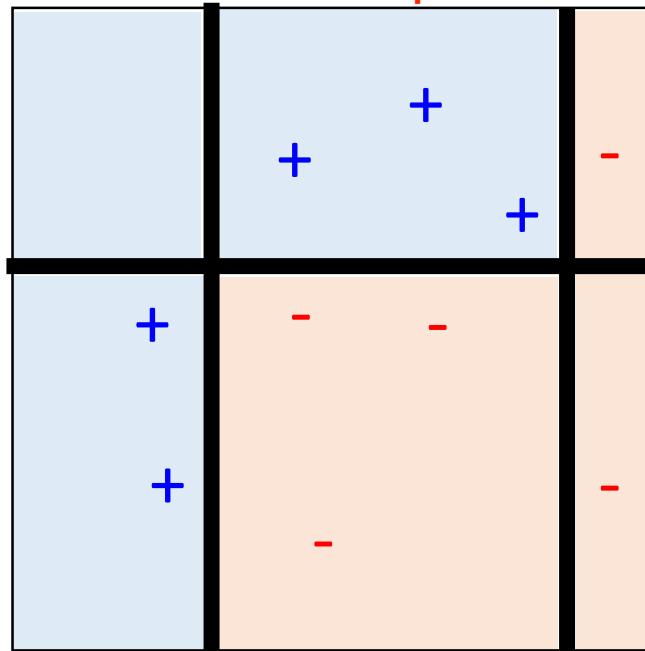
全部加起來看正負

- Final Classifier: $H(x) = \text{sign}(\sum_{t=1}^T \alpha_t f_t(x))$

分別乘上其weight (alpha)

$$\text{sign}(0.42 \begin{array}{|c|c|} \hline \text{+} & \text{+} \\ \hline \text{+} & \text{-} \\ \hline \end{array} + 0.66 \begin{array}{|c|c|} \hline \text{+} & \text{+} \\ \hline \text{+} & \text{-} \\ \hline \end{array} + 0.95 \begin{array}{|c|c|} \hline \text{+} & \text{+} \\ \hline \text{+} & \text{-} \\ \hline \end{array})$$

當組合起來這三個decision stump時可以得到正確率100的結果



Warning of Math

$$H(x) = \text{sign} \left(\sum_{t=1}^T \alpha_t f_t(x) \right) \quad \begin{array}{l} \text{epsilon t: classifier } f_t \text{ 的error rate} \\ \alpha_t = \ln \sqrt{(1 - \varepsilon_t) / \varepsilon_t} \end{array}$$

As we have more and more f_t (T increases), $H(x)$ achieves smaller and smaller error rate on training data.

隨著iteration越高adaboost結果越好

Error Rate of Final Classifier

- Final classifier: $H(x) = \text{sign}(\sum_{t=1}^T \alpha_t f_t(x))$
 $\alpha_t = \ln \sqrt{(1 - \varepsilon_t) / \varepsilon_t}$ $g(x)$

Training Data Error Rate

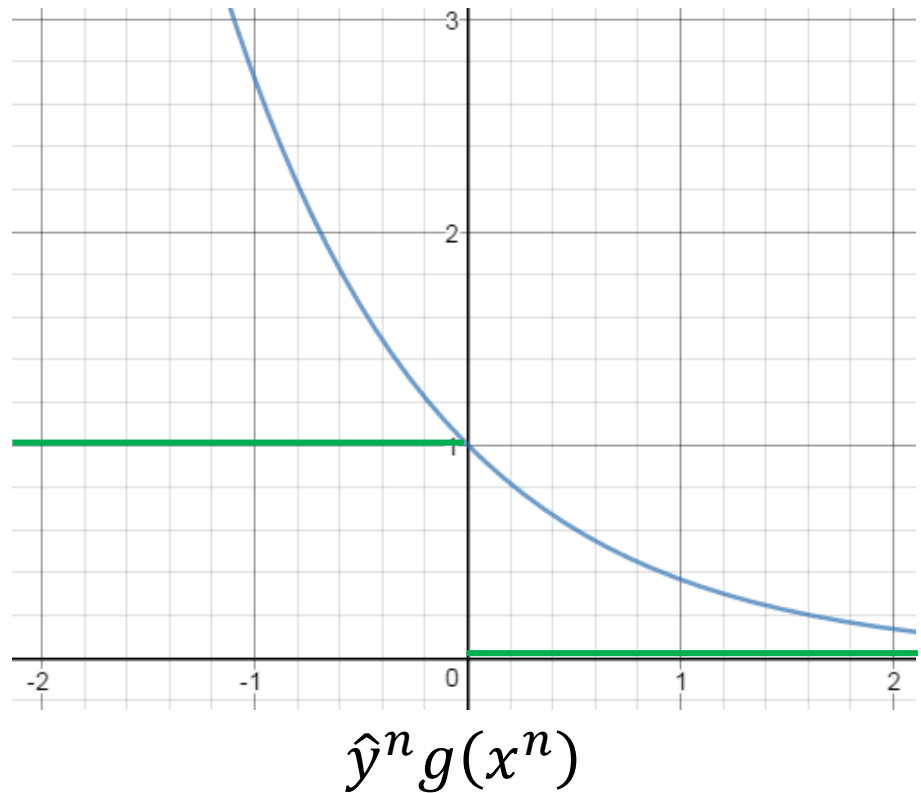
$$= \frac{1}{N} \sum_n \delta(H(x^n) \neq \hat{y}^n)$$

$$= \frac{1}{N} \sum_n \delta(\hat{y}^n g(x^n) < 0)$$

判斷同號異號

$$\leq \frac{1}{N} \sum_n \exp(-\hat{y}^n g(x^n))$$

upper bound



證明upper bound會越來越小
Training Data Error Rate

$$\leq \frac{1}{N} \sum_n \exp(-\hat{y}^n g(x^n)) = \frac{1}{N} Z_{T+1}$$

$$g(x) = \sum_{t=1}^T \alpha_t f_t(x)$$

$$\alpha_t = \ln \sqrt{(1 - \varepsilon_t) / \varepsilon_t}$$

Z_t : the summation of the weights of training data for training f_t

What is Z_{T+1} =? $Z_{T+1} = \sum_n u_{T+1}^n$

第1個iteration

$$\left. \begin{array}{l} u_1^n = 1 \\ u_{t+1}^n = u_t^n \times \exp(-\hat{y}^n f_t(x^n) \alpha_t) \end{array} \right\} \text{控制增加或是減少weight} \quad u_{T+1}^n = \prod_{t=1}^T \exp(-\hat{y}^n f_t(x^n) \alpha_t)$$

$$Z_{T+1} = \sum_n \prod_{t=1}^T \exp(-\hat{y}^n f_t(x^n) \alpha_t)$$

\hat{y} : label data固定的，可以提出

只要證明train data的weight之sum越來越小即可得證

$$= \sum_n \exp \left(-\hat{y}^n \sum_{t=1}^T f_t(x^n) \alpha_t \right)$$

Training Data Error Rate

$$\leq \frac{1}{N} \sum_n \exp(-\hat{y}^n g(x^n)) = \frac{1}{N} Z_{T+1}$$

$$g(x) = \sum_{t=1}^T \alpha_t f_t(x)$$
$$\alpha_t = \ln \sqrt{(1 - \epsilon_t) / \epsilon_t}$$

$$Z_1 = N \quad (\text{equal weights})$$

$$Z_t = \underbrace{Z_{t-1} \epsilon_t}_{\text{Misclassified portion in } Z_{t-1}} \exp(\alpha_t) + \underbrace{Z_{t-1} (1 - \epsilon_t)}_{\text{Correctly classified portion in } Z_{t-1}} \exp(-\alpha_t)$$

Misclassified portion in Z_{t-1}

Correctly classified portion in Z_{t-1}

$$= Z_{t-1} \epsilon_t \sqrt{(1 - \epsilon_t) / \epsilon_t} + Z_{t-1} (1 - \epsilon_t) \sqrt{\epsilon_t / (1 - \epsilon_t)}$$

$$= Z_{t-1} \times 2\sqrt{\epsilon_t (1 - \epsilon_t)}$$
$$Z_{T+1} = N \prod_{t=1}^T 2\sqrt{\epsilon_t (1 - \epsilon_t)}$$

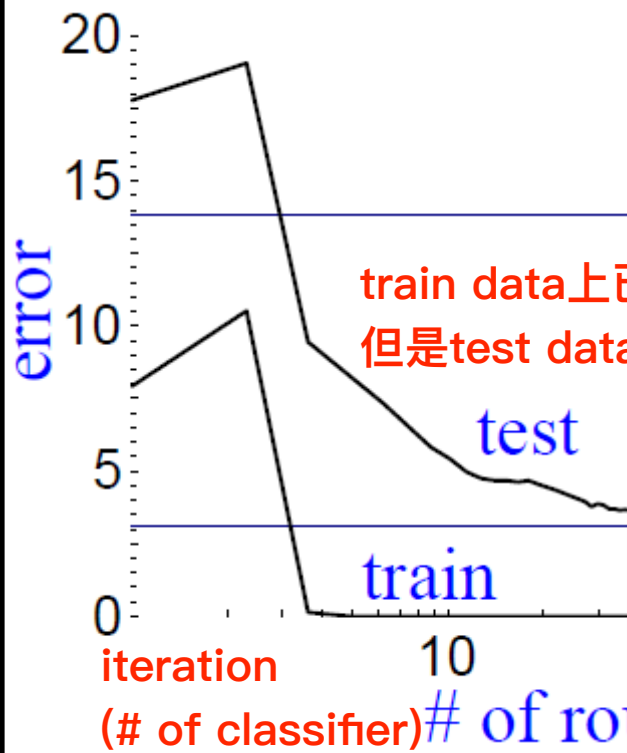
$$\text{Training Data Error Rate} \leq \prod_{t=1}^T \frac{2\sqrt{\epsilon_t (1 - \epsilon_t)}}{<1}$$

得證

Smaller and
smaller

End of Warning

Even though the training error is 0,
the testing error still decreases?



combined model

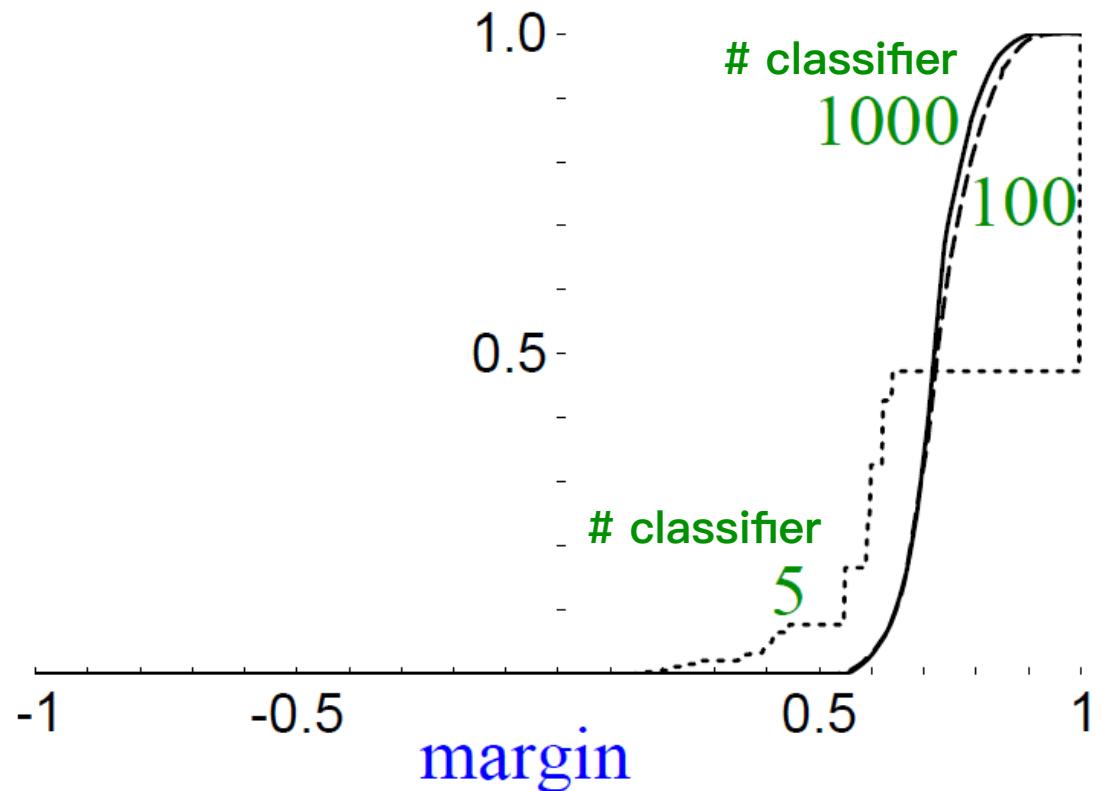
$$H(x) = \text{sign} \left(\sum_{t=1}^T \alpha_t f_t(x) \right)$$

$g(x)$

Margin = $\hat{y}g(x)$

與0的距離

cumulative distribution



為什麼可以增加margin? Large Margin?

$$H(x) = \text{sign} \left(\underbrace{\sum_{t=1}^T \alpha_t f_t(x)}_{g(x)} \right)$$

Training Data Error Rate =

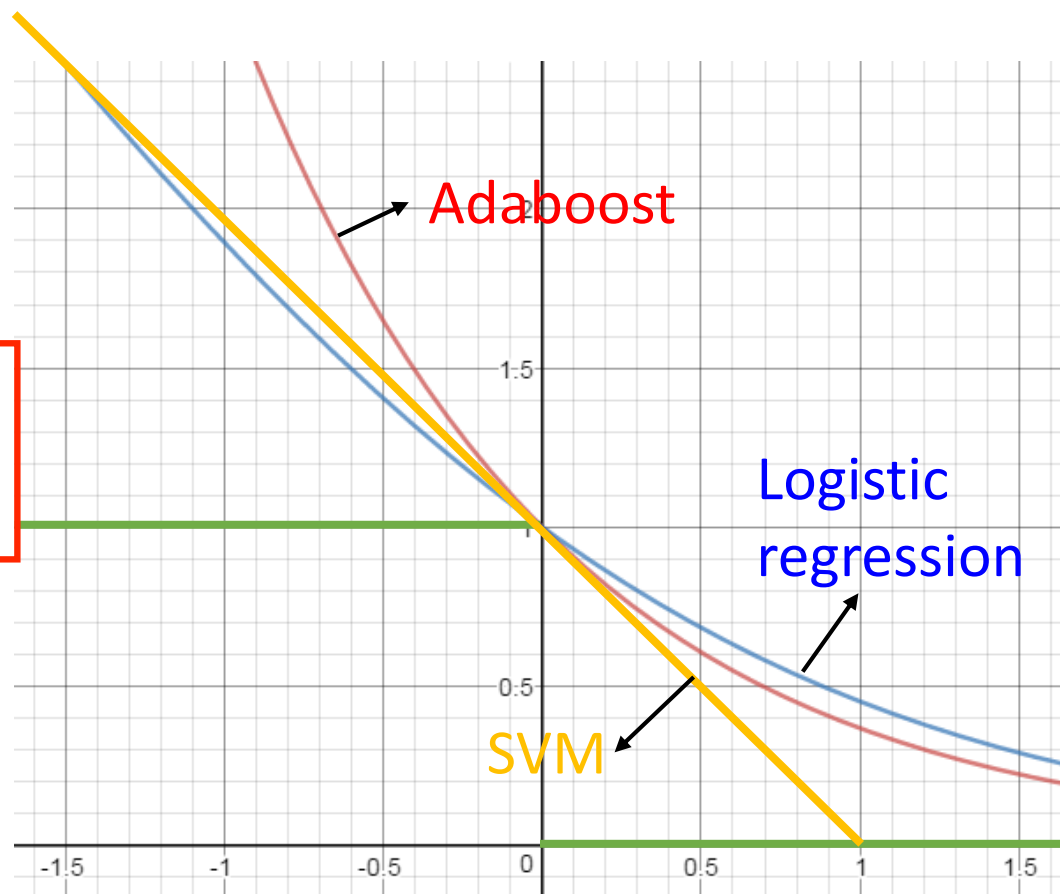
$$= \frac{1}{N} \sum_n \delta(H(x^n) \neq \hat{y}^n)$$

$$\leq \frac{1}{N} \sum_n \exp(-\hat{y}^n g(x^n))$$

adaboost

$$= \prod_{t=1}^T 2\sqrt{\epsilon_t(1-\epsilon_t)}$$

Getting smaller and
smaller as T increase

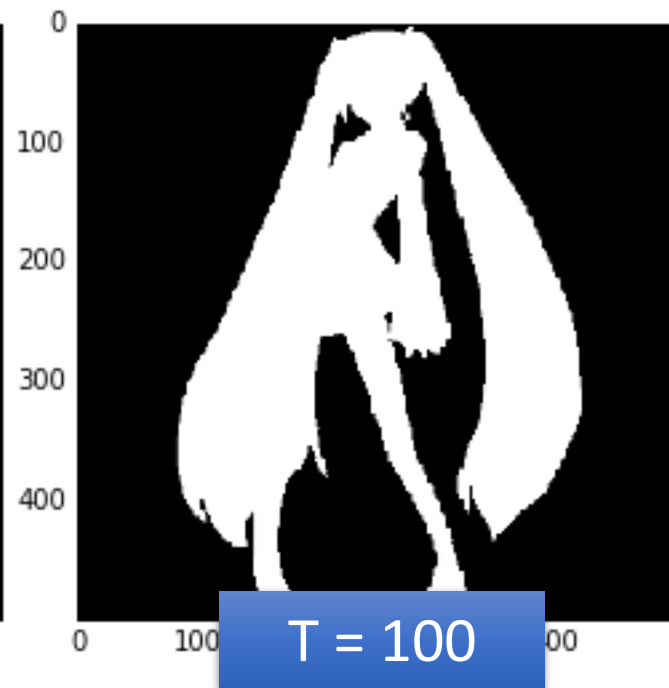
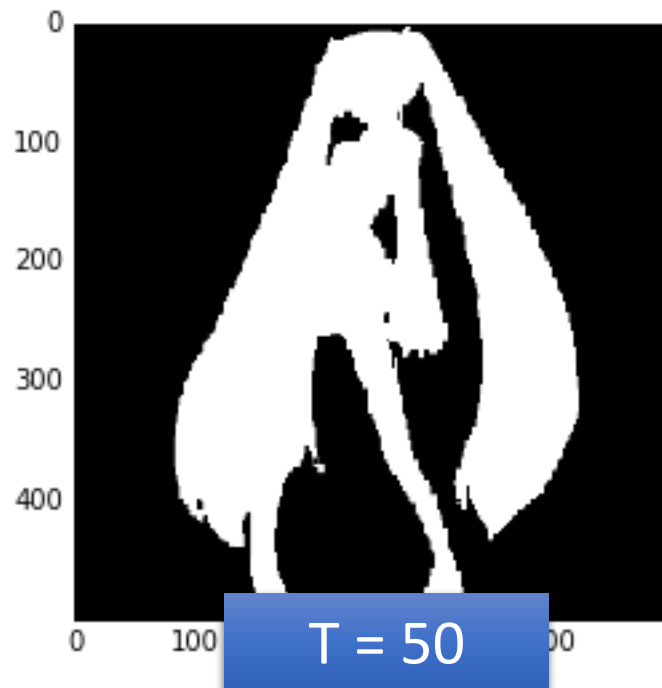
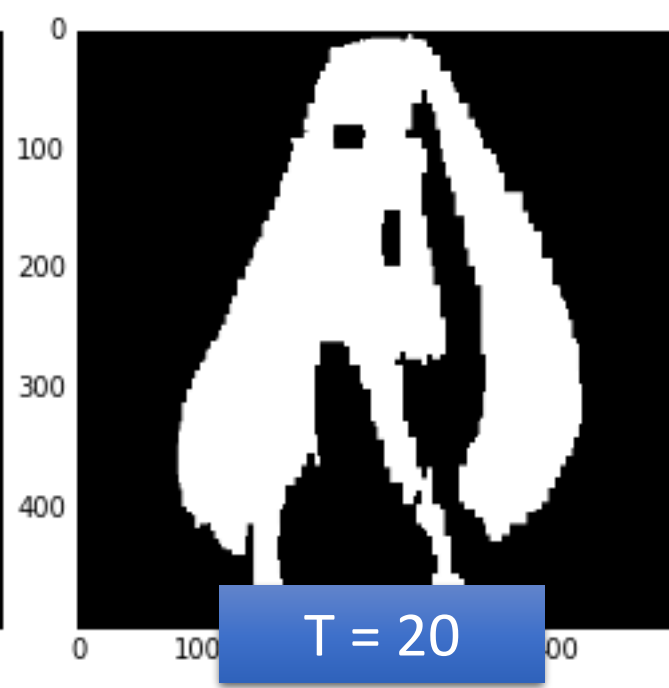
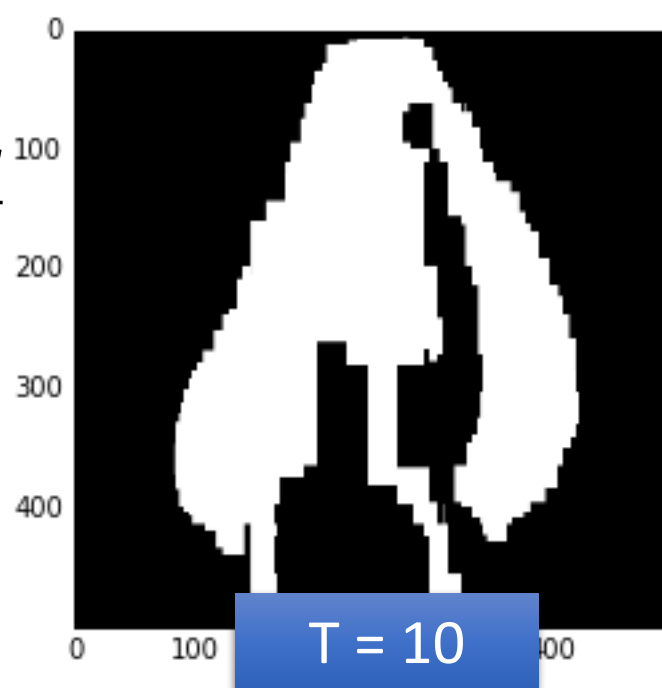


$\hat{y}^n g(x^n)$ 就算error rate=0，對於adaboost在往右的同時error可以在更小

Experiment:
Function of Miku

Adaboost
+Decision Tree

(depth = 5)



To learn more ...

- Introduction of Adaboost:
 - Freund; Schapire (1999). "A Short Introduction to Boosting"
- Multiclass/Regression
 - Y. Freund, R. Schapire, "A Decision-Theoretic Generalization of on-Line Learning and an Application to Boosting", 1995.
 - Robert E. Schapire and Yoram Singer. Improved boosting algorithms using confidence-rated predictions. In Proceedings of the Eleventh Annual Conference on Computational Learning Theory, pages 80–91, 1998.
- Gentle Boost
 - Schapire, Robert; Singer, Yoram (1999). "Improved Boosting Algorithms Using Confidence-rated Predictions".

General Formulation of Boosting

- Initial function $g_0(x) = 0$
- For $t = 1$ to T :
 - Find a function $f_t(x)$ and α_t to improve $g_{t-1}(x)$
 - $g_{t-1}(x) = \sum_{i=1}^{t-1} \alpha_i f_i(x)$
 - $g_t(x) = g_{t-1}(x) + \alpha_t f_t(x)$
- Output: $H(x) = \text{sign}(g_T(x))$

What is the learning target of $g(x)$? 如何找到 $f_t(x)$

Minimize $L(g) = \sum_n \underset{\text{loss function}}{l(\hat{y}^n, g(x^n))} = \sum_n \exp(-\hat{y}^n g(x^n))$ 將loss function定義如下

Gradient Boosting

- Find $g(x)$, minimize $L(g) = \sum_n \exp(-\hat{y}^n g(x^n))$
 - If we already have $g(x) = g_{t-1}(x)$, how to update $g(x)$?

想像：改變 g 在某一個點的值對 L 的影響有多大

把function G 對目標函式微分

Gradient Descent:

$$g_t(x) = g_{t-1}(x) - \eta \left. \frac{\partial L(g)}{\partial g(x)} \right|_{g(x) = g_{t-1}(x)}$$


Same direction

$$- \sum_n \exp(-\hat{y}^n g_{t-1}(x^n)) (-\hat{y}^n)$$

Boosting角度

$$g_t(x) = g_{t-1}(x) + \alpha_t f_t(x)$$

Gradient Boosting

vector $f_t(x)$  $\sum_n \exp(-\hat{y}^n g_t(x^n)) (\hat{y}^n)$ vector

Same direction

We want to find $f_t(x)$ maximizing

判斷正確
Minimize Error

值越大代表方向越一致

$$\sum_n \frac{\exp(-\hat{y}^n g_{t-1}(x^n))}{\text{example weight } u_t^n} (\hat{y}^n) f_t(x^n)$$

Same sign

$$u_t^n = \exp(-\hat{y}^n g_{t-1}(x^n)) = \exp\left(-\hat{y}^n \sum_{i=1}^{t-1} \alpha_i f_i(x^n)\right)$$

$$= \prod_{i=1}^{t-1} \exp(-\hat{y}^n \alpha_i f_i(x^n))$$

Exactly the weights we obtain in Adaboost

Gradient Boosting

先找出 f_t 後固定住，窮舉所有 α (learning rate) 讓loss最小

- Find $g(x)$, minimize $L(g) = \sum_n \exp(-\hat{y}^n g(x^n))$

$$g_t(x) = g_{t-1}(x) + \alpha_t f_t(x)$$

α_t is something like learning rate

Find α_t minimizing $L(g_{t+1})$

$$\begin{aligned} L(g) &= \sum_n \exp(-\hat{y}^n (g_{t-1}(x) + \alpha_t f_t(x))) \\ &= \sum_n \exp(-\hat{y}^n g_{t-1}(x)) \exp(-\hat{y}^n \alpha_t f_t(x)) \\ &= \sum_{\hat{y}^n \neq f_t(x)} \exp(-\hat{y}^n g_{t-1}(x^n)) \exp(\alpha_t) \\ &\quad + \sum_{\hat{y}^n = f_t(x)} \exp(-\hat{y}^n g_{t-1}(x^n)) \exp(-\alpha_t) \end{aligned}$$

Find α_t
such that

$$\frac{\partial L(g)}{\partial \alpha_t} = 0$$

$$\alpha_t = \ln \sqrt{(1 - \varepsilon_t) / \varepsilon_t}$$

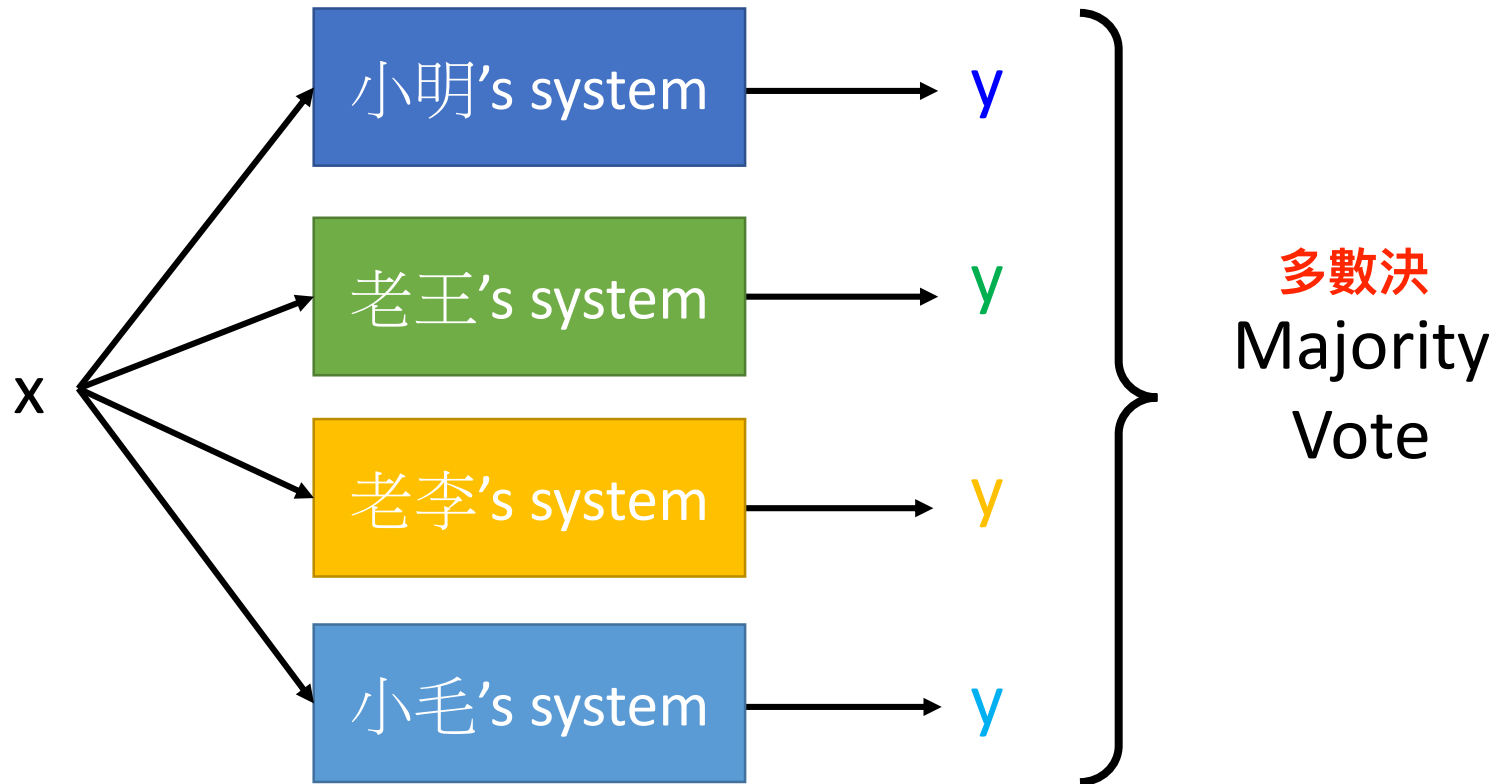
Adaboost!

Cool Demo

- http://arogozhnikov.github.io/2016/07/05/gradient_boosting_playground.html

Ensemble: Stacking

Voting



然而這樣做需要將training data分兩群，一部份train各自的model，一部分拿來做final classifier的training，避免有些異常的model雖然accuracy=100%但是是因為overfitting，這樣final classifier只會參考那個fit 100%的壞model

Stacking

