# Classification:
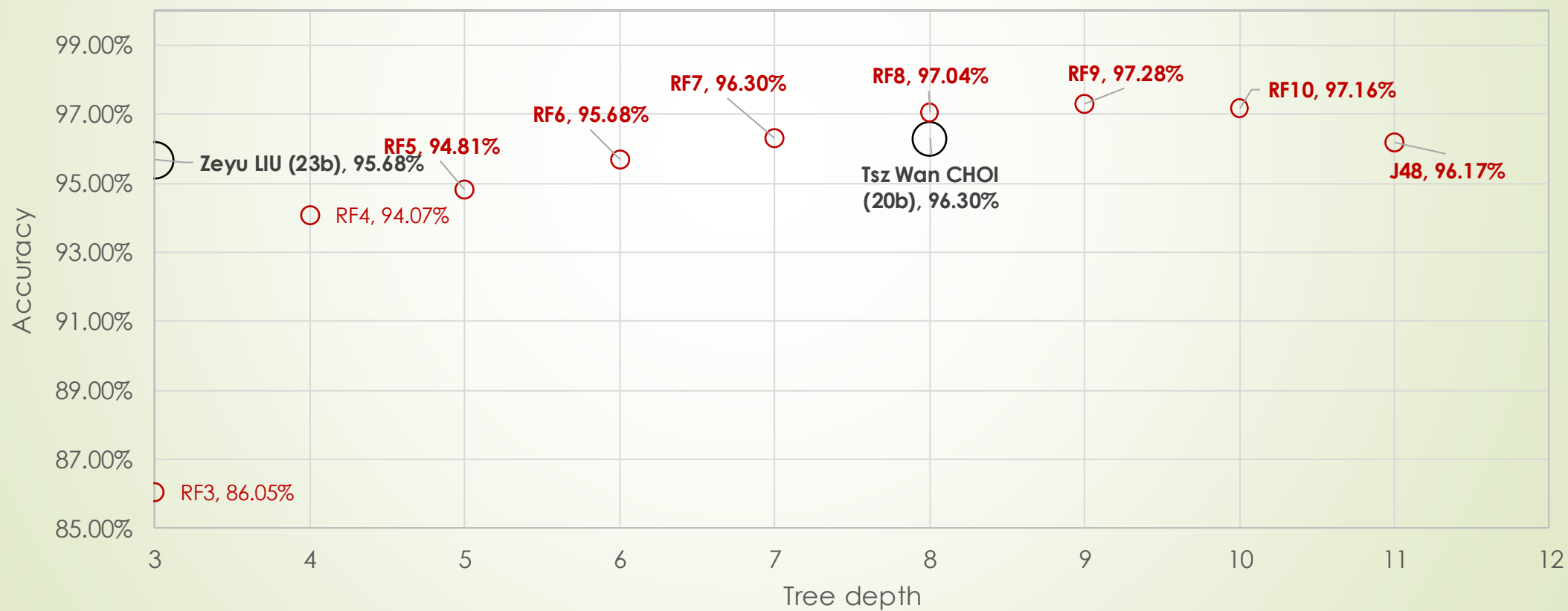# Ensemble Methods

CS5483 Data Warehousing and Data Mining
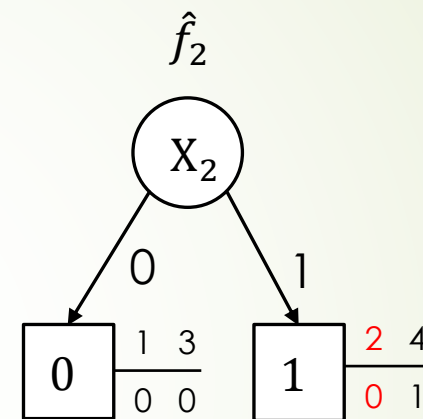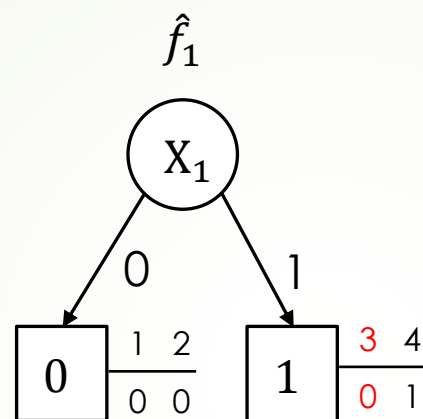
# Man vs Machine Rematch

# Segment Challenge Results
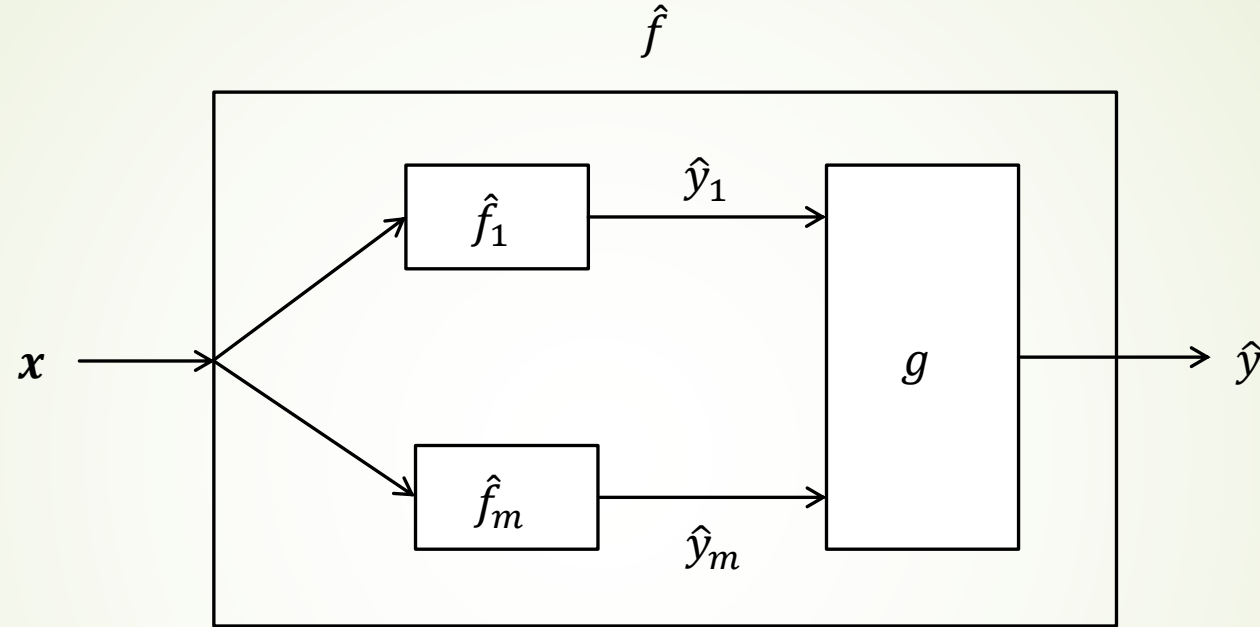
# Two heads are better than one

- [Bing](#)/[Baidu](#)/[Google](#) translation.
- The story in [Chinese](#) and its translation to [English](#).
- Can we combine two poor classifiers into a good classifier?
- What is the benefit of doing so?

$$\hat{f}_1 \qquad\qquad \hat{f}_2$$

|   | $X_1$ | $X_2$ | $Y$ |
|---|---|---|---|
| 1. | 0 | 0 | 0 |
| 2. | 0 | 1 | 0 |
| 3. | 1 | 0 | 0 |
| 4. | 1 | 1 | 1 |

$\hat{f}_1$: node $X_1$; branch 0 → leaf $0$ ($\frac{1\ 2}{0\ 0}$); branch 1 → leaf $1$ ($\frac{3\ 4}{0\ 1}$)

$\hat{f}_2$: node $X_2$; branch 0 → leaf $0$ ($\frac{1\ 3}{0\ 0}$); branch 1 → leaf $1$ ($\frac{2\ 4}{0\ 1}$)

- Accuracies of $\hat{f}_1$ and $\hat{f}_2$ are both _____%. Are they good?

- Can we combine them into a better classifier $\hat{f}(x) := g\left(\hat{f}_1(x), \hat{f}_2(x)\right)$?

- _____$\{\hat{f}_1(x), \hat{f}_2(x)\}$ achieves an accuracy of _____%.
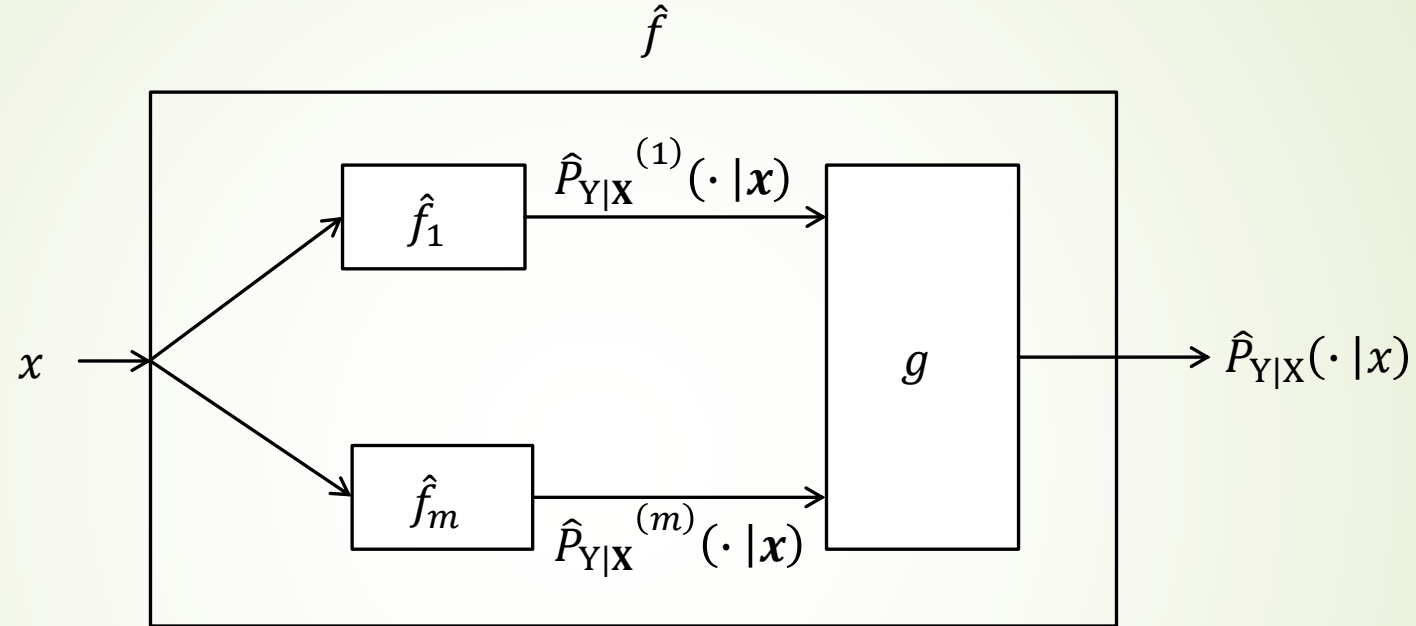
- How does it work in general?

# Architecture



1. **Base classifiers** $\hat{f}_j$'s are simple but possibly have weak preliminary predictions $\hat{y}_j$'s.
2. **Combined classifier** $\hat{f}$ uses the **combination rule** $g$ to merge $\hat{y}_j$'s into a good final prediction $\hat{y}$.

# Architecture for probabilistic classifiers

$$\hat{f}$$



1. **Base classifiers** $\hat{f}_j$'s are simple but possibly have weak probability estimates $\hat{P}_{Y|\mathbf{X}}^{(j)}(\cdot\,|x)$.

2. **Combined classifier** $\hat{f}$ uses the **combination rule** $g$ to merge $\hat{P}_{Y|\mathbf{X}}^{(j)}(\cdot\,|\boldsymbol{x})$'s into a good final prediction $\hat{P}_{Y|\mathbf{X}}(\cdot\,|x)$.

# How to get good performance?

- Reduce **risk** by avoiding *underfitting* and *overfitting*.
- For many loss functions $L$ (0-1 loss, sum of squared error, ...):

$$\overbrace{E\left[L\left(Y, f_{\mathbf{W}}(X)\right)\right]}^{\text{Risk}} \leq \overbrace{E\left[L\left(Y, \bar{f}(X)\right)\right]}^{\text{Bias}} + \overbrace{E\left[L\left(\bar{f}(X), f_{\mathbf{W}}(X)\right)\right]}^{\text{Variance}}$$

where

- $\bar{f} := x \mapsto E[f_{\mathbf{W}}(x)]$ is the **expected predictor;** (W is a random variable. Why?)
- **Variance** is the dependence of $f_{\mathbf{W}}(X)$ on the data aka <u>overfitting/underfitting</u>; and
- **Bias** is the deviation of $\bar{f}(X)$ from Y aka <u>overfitting/underfitting</u>.
- See <u>Bias-variance trade-off</u>.

# Bias and variance for probabilistic classifiers

 For probabilistic classifiers,

$$\overbrace{E\left[L\left(P_{Y|X}(\cdot\,|X), P_{\widehat{Y}|X,W}(\cdot\,|X,W)\right)\right]}^{\text{Risk}} \leq \overbrace{E\left[L\left(P_{Y|X}(\cdot\,|X), P_{\widehat{Y}|X}(\cdot\,|X)\right)\right]}^{\text{Bias}} + \overbrace{I(\widehat{Y};W|X)}^{\text{Variance}}$$

where

 $f_w(x) := P_{\widehat{Y}|X,W}(\cdot\,|x,w)$ implies $\bar{f}(x) = E\left[P_{\widehat{Y}|X,W}(\cdot\,|x,W)\right] = P_{\widehat{Y}|X}(\cdot\,|x)$, called m_____;

 $P_{Y|X}(\cdot\,|X)$ instead of $Y$ is used as the ground truth;

 information (or Kullback-Leibler) divergence is used as the loss function

$$L(Q,P) := D_{\mathrm{KL}}(P\|Q) := \int_{\mathcal{Y}} (dP) \log \frac{dP}{dQ};\ \ \text{and}$$
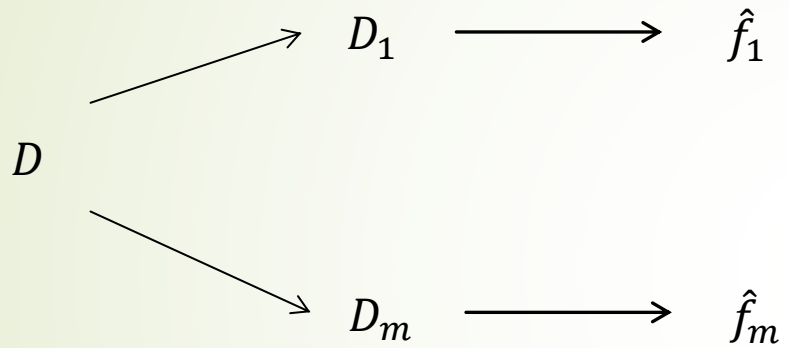
 variance becomes the mutual information

$$E\left[D_{\mathrm{KL}}\left(P_{\widehat{Y}|X,W}(\cdot\,|X,W)\|P_{\widehat{Y}|X}(\cdot\,|X)\right)\right] = I(\widehat{Y};W|X) \qquad \because I(X;W) = 0.$$

# How to reduce variance and bias?

- Base classifiers should be **d**_____, i.e., capture **as many different pieces of relevant information** as possible to reduce _____.

- The combination rule should reduce _____ by **smoothing out the noise** while **aggregating relevant information** into the final decision.
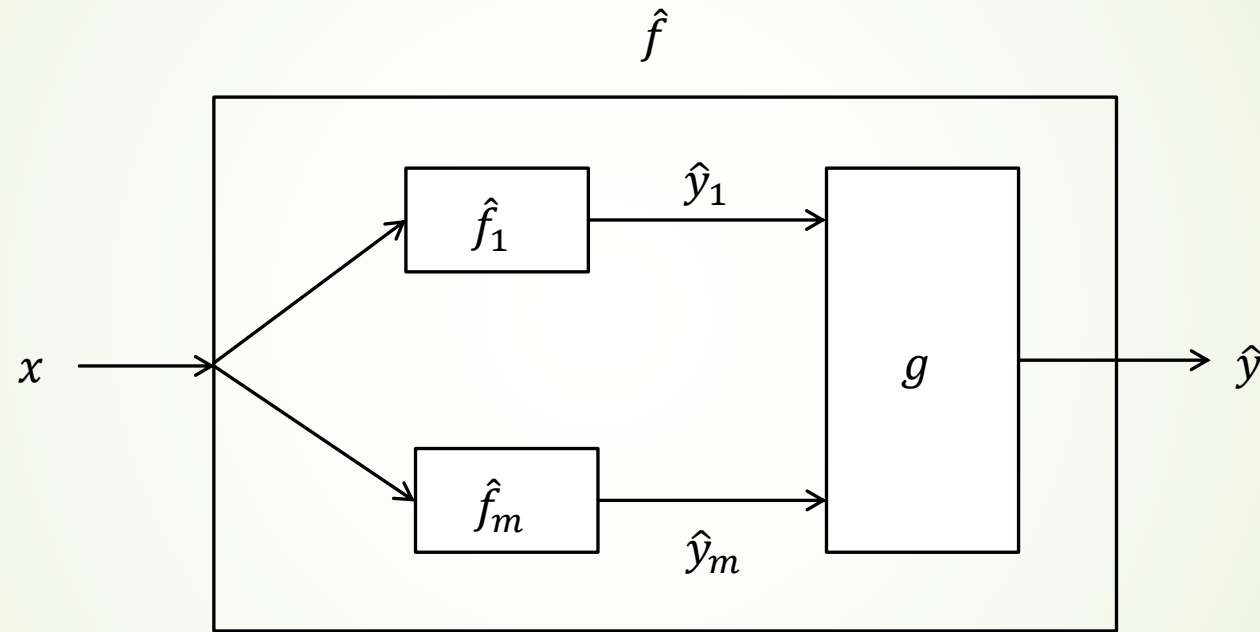
# Bagging (**B**ootstrap **Ag**gregation)
## Base classifiers

$$D \nearrow D_1 \longrightarrow \hat{f}_1$$

$$D \searrow D_m \longrightarrow \hat{f}_m$$

- Construct $m$ bootstrap samples.
- Construct a base classifier for each bootstrap sample.

# Bagging (**B**ootstrap **Ag**gregation)
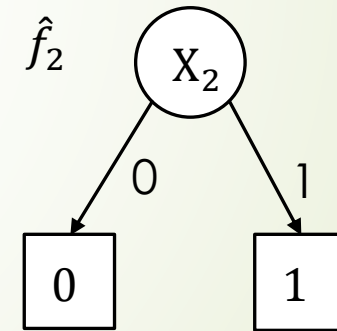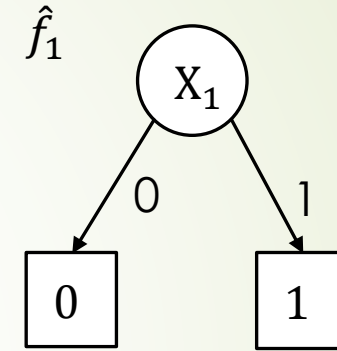## *Majority voting*

$$\hat{f}$$



$$\hat{f}(x) := \arg\max_{\hat{y}} \overbrace{\sum_{j} \mathbb{1}\left(\hat{f}_j(x) = \hat{y}\right)}^{|\{j|\hat{f}_j(x)=\hat{y}\}|=}$$

# Example

|    | $X_1$ | $X_2$ | Y |
|----|-------|-------|---|
| 1. | 0     | 0     | 0 |
| 2. | 0     | 1     | 0 |
| 3. | 1     | 0     | 0 |
| 4. | 1     | 1     | 1 |

|    | $X_1$ | $X_2$ | Y |
|----|-------|-------|---|
| 1. | 0     | 0     | 0 |
| 2. | 0     | 1     | 0 |
| 2. | 0     | 1     | 0 |
| 4. | 1     | 1     | 1 |

|    | $X_1$ | $X_2$ | Y |
|----|-------|-------|---|
| 1. | 0     | 0     | 0 |
| 3. | 1     | 0     | 0 |
| 3. | 1     | 0     | 0 |
| 4. | 1     | 1     | 1 |

$\hat{f}_1$



$\hat{f}_2$

| | $X_1$ | $X_2$ | $Y$ | $\hat{f}_1$ | $\hat{f}_2$ | $\hat{f}$ |
|---|---|---|---|---|---|---|
| 1. | 0 | 0 | 0 | 0 | 0 | 0 |
| 2. | 0 | 1 | 0 | 0 | 1 | ? |
| 3. | 1 | 0 | 0 | 1 | 0 | ? |
| 4. | 1 | 1 | 1 | 1 | 1 | 1 |

0

1

with equal probability

Accuracy = _____%

$\hat{f}_1$

$X_1$

0        1

0        1

$\hat{f}_2$

$X_2$

0        1

0        1

# Is it always good to follow the majority?

$\hat{f}_1$



| | $X_1$ | $X_2$ | Y | $\hat{f}_1$ | $\hat{f}_2$ | $\hat{f}$ | |
|---|---|---|---|---|---|---|---|
| 1. | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 2. | 0 | 1 | 0 | 0 | 1 | ? | |
| 3. | 1 | 0 | 0 | 1 | 0 | ? | |
| 4. | 1 | 1 | 1 | 1 | 1 | 1 | |

Accuracy = _____%

$\hat{f}_2$



- It is beneficial to return 0 more often because _____.
- How to do this in general?

# Sum rule and threshold moving

- $\hat{f}(x) = 1$ iff

$$\frac{1}{2}\left[\hat{f}_1(x) + \hat{f}_2(x)\right] > \underline{\hspace{3cm}}$$

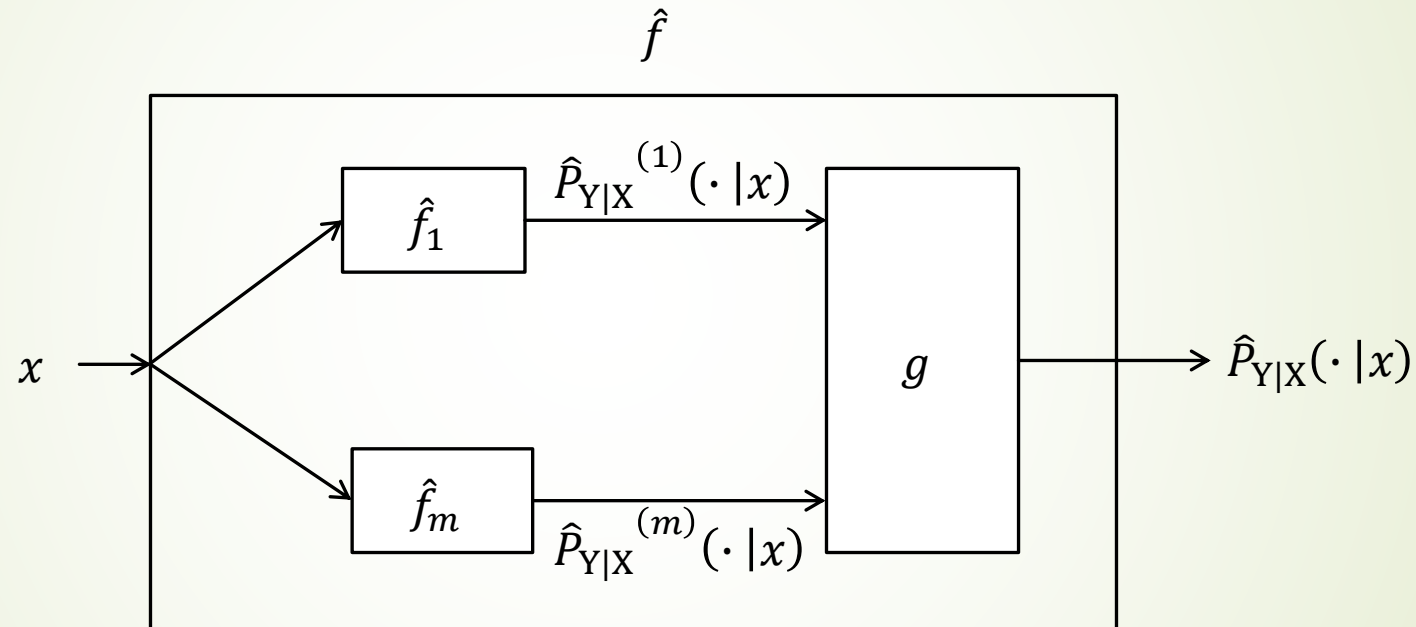- **Binary classification:** Choose $\hat{f}(x) = 1$ iff

$$\frac{1}{m}\sum_t \hat{f}_t(x) > \gamma$$

  for some chosen threshold $\gamma$.

- What about multi-class classification?

# Bagging (**B**ootstrap **Ag**gregation)
## Average of probabilities



$$\hat{f}(\boldsymbol{x}) := \frac{1}{m}\sum_t \widehat{f_t}(\boldsymbol{x})$$

# Other techniques to diversify base classifiers

- **Random forest**: Bagging with modified decision tree induction
  - **Forest-RI**: For each split, consider **random i**_____ **s**_____ where only $F$ randomly chosen features are considered.
  - **Forest-RC**: For each split, consider $F$ **random l**_____ **c**_____ of $L$ randomly chosen features.
- **Voting** (weka.classifier.meta.vote) and **Stacking** (weka.classifier.meta.stacking):
  - Use different classification algorithms.
- **Adaptive boosting (Adaboost)**:
  - Each base classifier tries to _____ made by previous base classifiers.

# Other techniques to combine decisions

- **Random forest**
  - Majority voting
  - Average of probabilities
- **Voting**
  - Majority voting or median
  - Average/product/minimum/maximum probabilities
- **Stacking**
  - Use a meta classifier.
- **Adaptive boosting (Adaboost)** - 2003 Gödel Prize winner
  - Weighted majority voting

# What is Adaboost?

- An ensemble method that **learns from mistakes**:
  - Combined classifier:
    - Majority voting but with more weight on more accurate base classifier.

    $$\hat{f}(x) := \arg\max_{\hat{y}} \sum_t w_t \cdot \mathbb{1}\big(\widehat{f_t}(x) = \hat{y}\big)$$

    where $w_t := \frac{1}{2}\ln\frac{1-\text{error}(\hat{f}_t)}{\text{error}(\hat{f}_t)}$ is the amount of say of $\hat{f}_t$ and
    $\text{error}(\hat{f}_t)$ is the error rate w.r.t. $D_t$. (See the precise formula below.)
  - Base classifiers:
    - Train $\widehat{f_t}$ sequentially in $t$ on $D_t$ obtained by
    - Bagging $(x_i, y_i) \in D$ with

    $$p_i^{(t)} := \frac{p_i^{(t-1)}}{Z_t} \times \begin{cases} e^{w_{t-1}}, & \hat{f}_{t-1}(x_i) \neq y_i \text{ (incorrectly classified example)} \\ e^{-w_{t-1}}, & \text{otherwise (correctly classified example).} \end{cases}$$
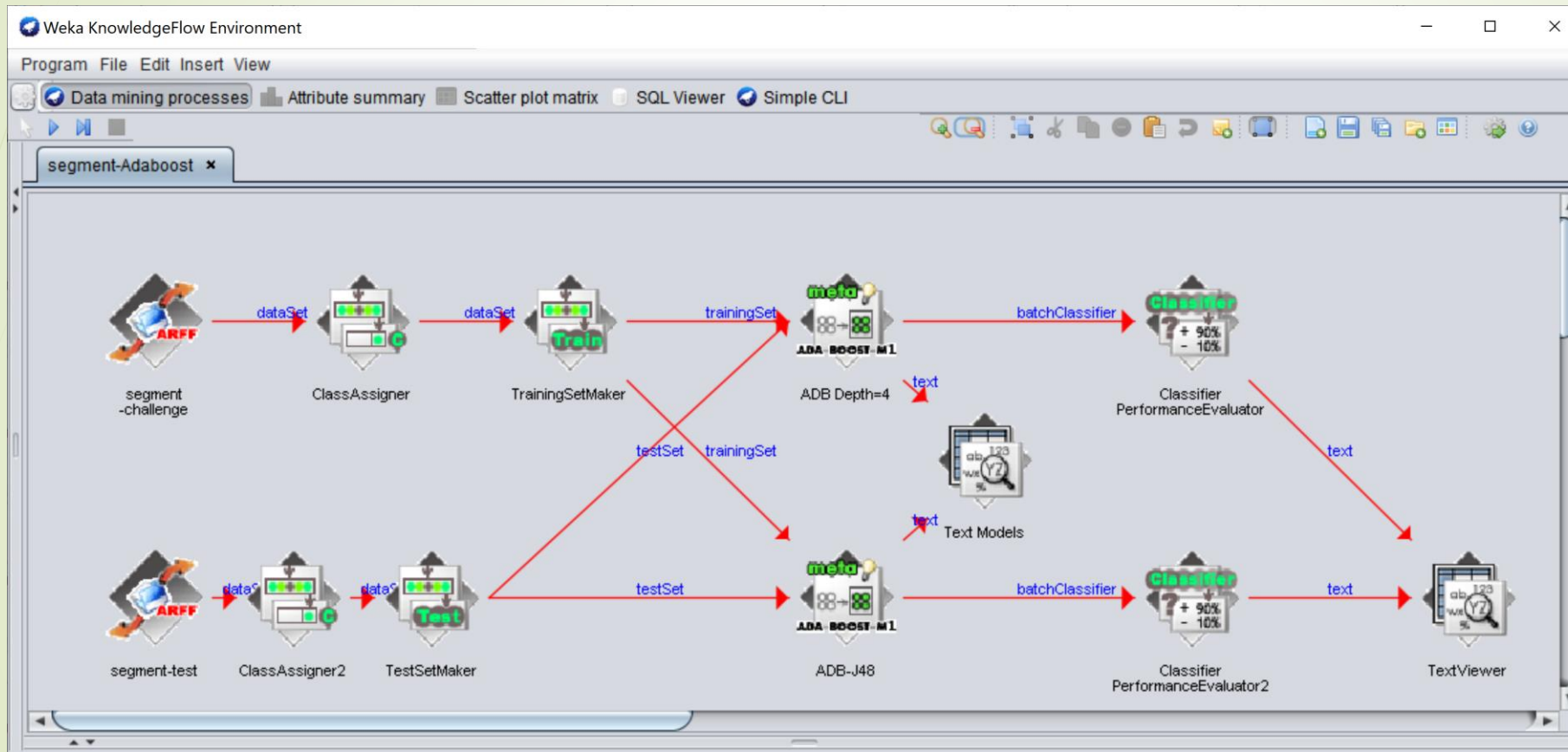
    starting with $p_i^{(1)} := \frac{1}{|D|}$ and with $Z_t > 0$ chosen so that $\sum_i p_i^{(t)} = 1$.
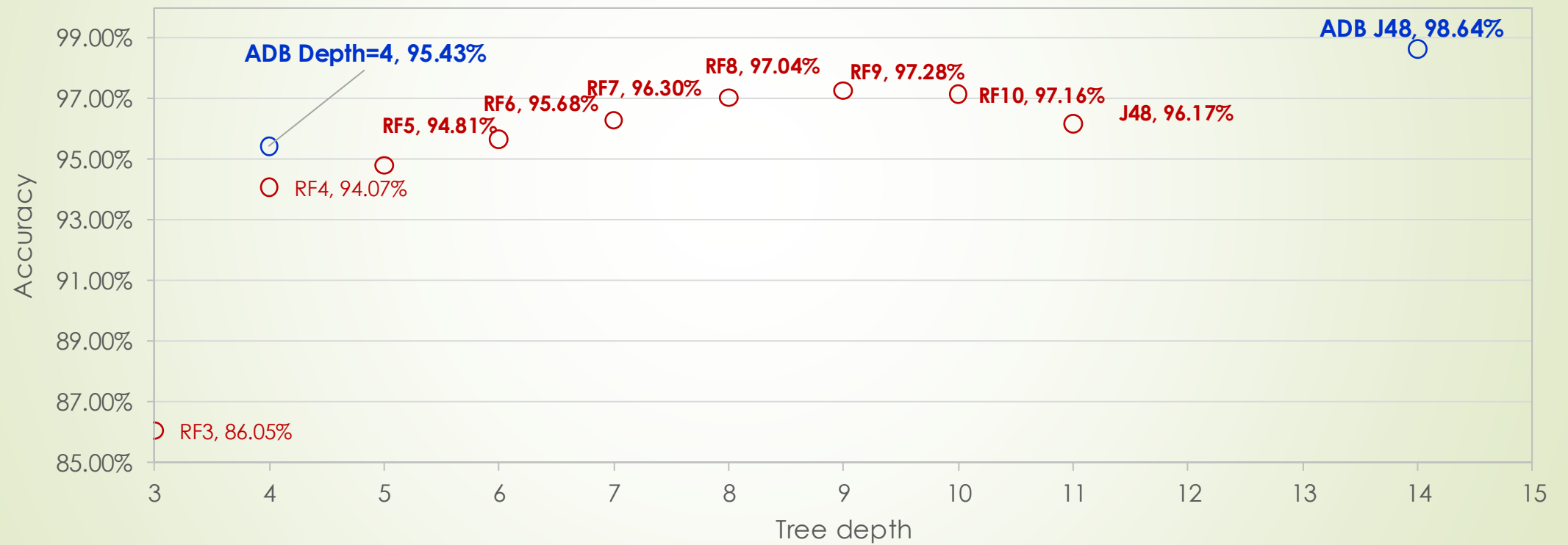    - Compute the error rate

    $$\text{error}(\hat{f}_t) := \sum_i p_i^{(t)} \cdot \mathbb{1}\big(\widehat{f_t}(x_i) \neq y_i\big)$$

# Machine vs Machine

# Machine vs Machine



ADB J48, 98.64%

ADB Depth=4, 95.43%

RF8, 97.04%

RF9, 97.28%

RF7, 96.30%

RF6, 95.68%

RF10, 97.16%

RF5, 94.81%

J48, 96.17%

RF4, 94.07%

RF3, 86.05%

Accuracy

Tree depth

# References

- 8.6 Techniques to improve classification accuracy

- [Witten11] Chapter 8

- *Optional:*

  - Breiman, L. (1996). "Bagging predictors." *Machine learning*, 24(2), 123-140.

  - Breiman, L. (2001). "Random forests." *Machine learning*, 45(1), 5-32.

  - Freund Y, Schapire R, Abe N. "A short introduction to boosting." Journal-Japanese Society For Artificial Intelligence. 1999 Sep 1;14(771-780):1612.

  - Zhu, H. Zou, S. Rosset, T. Hastie, "Multi-class AdaBoost", 2009.