# Lecture 23: Gradient Boosting, Bagging, Decision Forest CS480/680 Intro to Machine Learning

2023-4-4

Pascal Poupart
David R. Cheriton School of Computer Science

UNIVERSITY OF
**WATERLOO**

# Gradient Boosting

- AdaBoost designed for classification

- How can we use boosting for regression?

- Answer: **Gradient Boosting**

UNIVERSITY OF
**WATERLOO**

# Gradient Boosting

Idea:

- Predictor $f_k$ at stage $k$ incurs loss $L(f_k(\boldsymbol{x}), y)$

- Train $h_{k+1}$ to approximate negative gradient:

$$h_{k+1}(\boldsymbol{x}) \approx -\frac{\partial L(f_k(\boldsymbol{x}), y)}{\partial f_k(\boldsymbol{x})}$$

- Update predictor by adding a multiple $\eta_{k+1}$ of $h_{k+1}$:

$$f_{k+1}(\boldsymbol{x}) \leftarrow f_k(\boldsymbol{x}) + \eta_{k+1}\, h_{k+1}(\boldsymbol{x})$$

UNIVERSITY OF
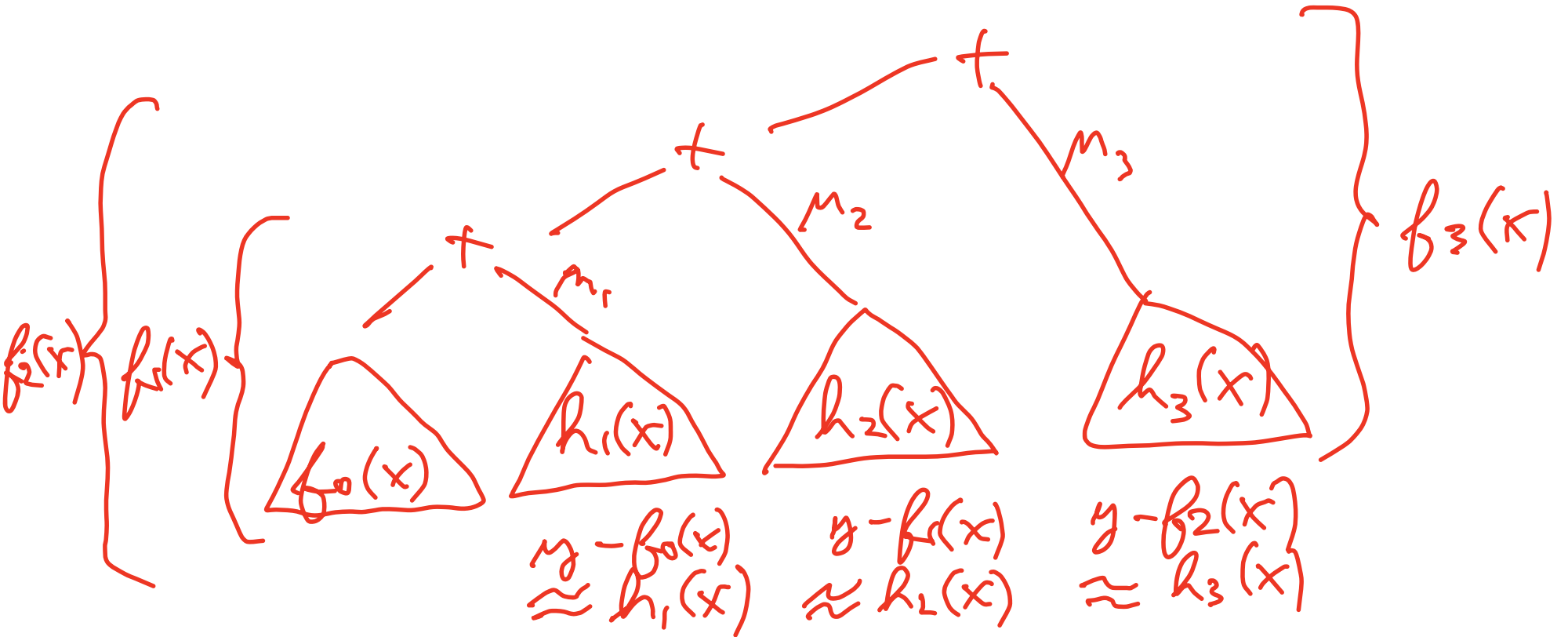WATERLOO

# Squared Loss

- Consider **squared loss**

$$L(f_k(\boldsymbol{x}_n), y_n) = \frac{1}{2}(f_k(\boldsymbol{x}_n) - y_n)^2$$

- Negative gradient corresponds to **residual** $r$

$$-\frac{\partial L(f_k(\boldsymbol{x}_n), y_n)}{\partial f_k(\boldsymbol{x}_n)} = y_n - f_k(\boldsymbol{x}_n) = r_n$$

- Train **base learner** $h_{k+1}$
  with **residual dataset** $\{(\boldsymbol{x}_n, r_n)_{\forall n}\}$

- Base learner $h_{k+1}$ can be any **non-linear predictor** (often a small decision tree)

UNIVERSITY OF
**WATERLOO**

# Illustration



$f_3(x)$, $f_2(x)$, $f_1(x)$, $M_1$, $M_2$, $M_3$

$h_0(x)$, $h_1(x)$, $h_2(x)$, $h_3(x)$

$y - f_0(x) \approx h_1(x)$

$y - f_1(x) \approx h_2(x)$

$y - f_2(x) \approx h_3(x)$

UNIVERSITY OF
WATERLOO

# Gradient Boosting Algorithm

- Initialize predictor with a constant $c$:
$$f_0(\boldsymbol{x}_n) = argmin_c \sum_n L(c, y_n)$$

- For $k = 1$ to $K$ do

  - Compute pseudo residuals: $r_n = -\dfrac{\partial L(f_{k-1}(\boldsymbol{x}_n), y_n)}{\partial f_{k-1}(\boldsymbol{x}_n)}$

  - Train a base learner $h_k$ with residual dataset $\{(\boldsymbol{x}_n, r_n)_{\forall n}\}$

  - Optimize step length:
  $$\eta_k = argmin_\eta \sum_n L(f_{k-1}(\boldsymbol{x}_n) + \eta h_k(\boldsymbol{x}_n), y_n)$$

  - Update predictor: $f_k(\boldsymbol{x}) \leftarrow f_{k-1}(\boldsymbol{x}) + \eta_k h_k(\boldsymbol{x})$

UNIVERSITY OF
WATERLOO

# XGBoost

- eXtreme Gradient Boosting

  - Package optimized for speed and accuracy

  - XGBoost used in >12 winning entries for various challenges
    https://github.com/dmlc/xgboost/tree/master/demo#machine-learning-challenge-winning-solutions

UNIVERSITY OF
WATERLOO

# Boosting vs Bagging

- Review

Bagging

- majority vote
- Assumptions
  - independent hypotheses
  - hypotheses with similar accuracies

Boosting

- weighted predictions
- Allows:
  - Correlated hypotheses
  - Hypotheses with imbalanced accuracies
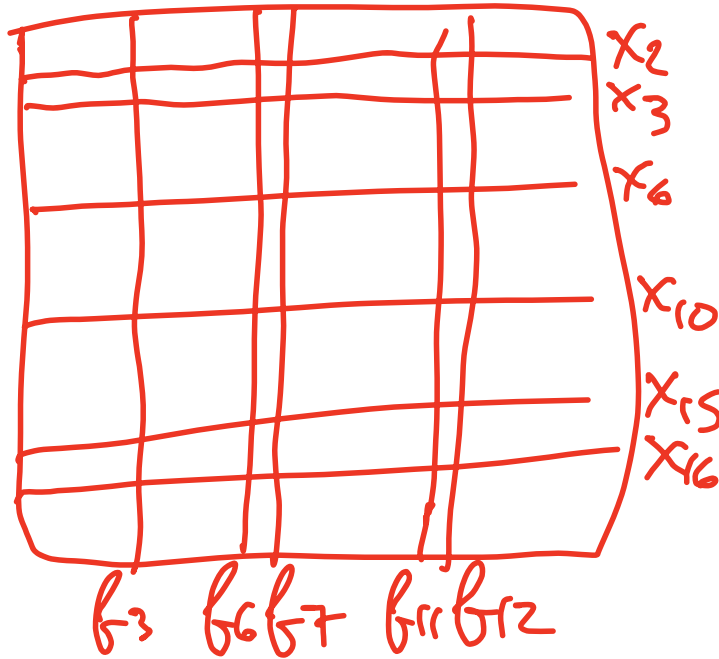
UNIVERSITY OF
WATERLOO

# Independent classifiers/predictors

- How can we obtain independent classifiers/predictors for bagging?

- Bootstrap sampling
  - Sample (without replacement) subset of data

- Random projection
  - Sample (without replacement) subset of features

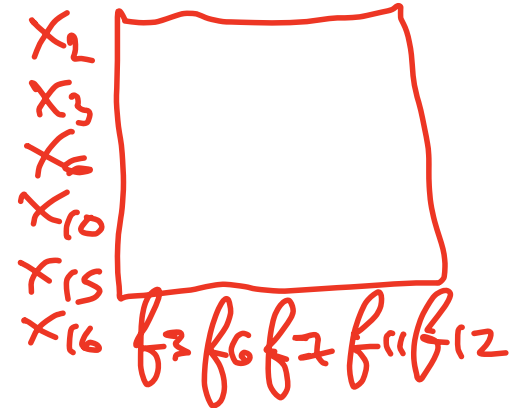- Learn different classifiers/predictors based on each data subset and feature subset

UNIVERSITY OF
**WATERLOO**

# Illustration of Bootstrap Sampling and Random Projection

Dataset

features

instance

$X_2$
$X_3$
$X_6$
$X_{10}$
$X_{15}$
$X_{16}$

$f_3 \; f_6 f_7 \; f_{11} f_{12}$

Sampled dataset

$X_2$
$X_3$
$X_6$
$X_{10}$
$X_{15}$
$X_{16}$ $f_3 \, f_6 \, f_7 \, f_{11} f_{12}$

UNIVERSITY OF
WATERLOO

# Bagging

For k = 1 to K

$\quad\quad D_k \leftarrow$ sample data subset

$\quad\quad F_k \leftarrow$ sample feature subset

$\quad\quad h_k \leftarrow$ train classifier/predictor based on $D_k$ and $F_k$

Classification: $majority(h_1(\boldsymbol{x}), \dots, h_K(\boldsymbol{x}))$

Regression: $average(h_1(\boldsymbol{x}), \dots, h_K(\boldsymbol{x}))$

**Random forest:** bag of decision trees

UNIVERSITY OF
WATERLOO

# Application: Xbox 360 Kinect

- Microsoft Cambridge

- Body part recognition: supervised learning

# Depth camera



- Kinect

Infrared image

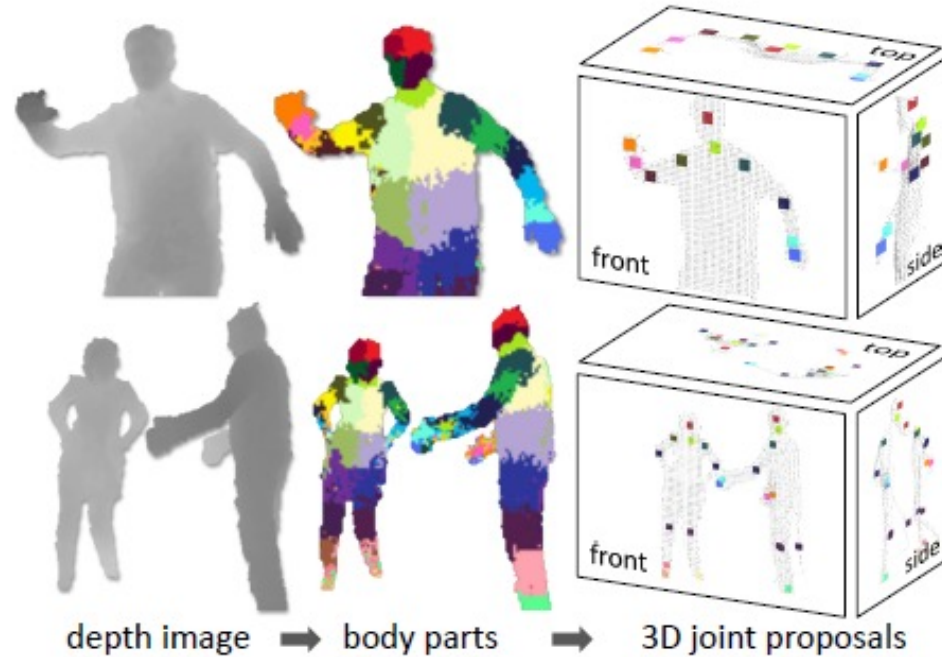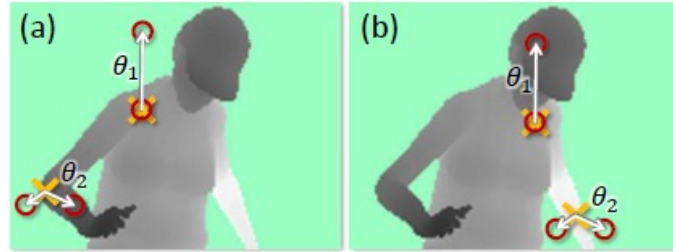Gray scale depth map

# Kinect Body Part Recognition

- Problem: label each pixel with a body part



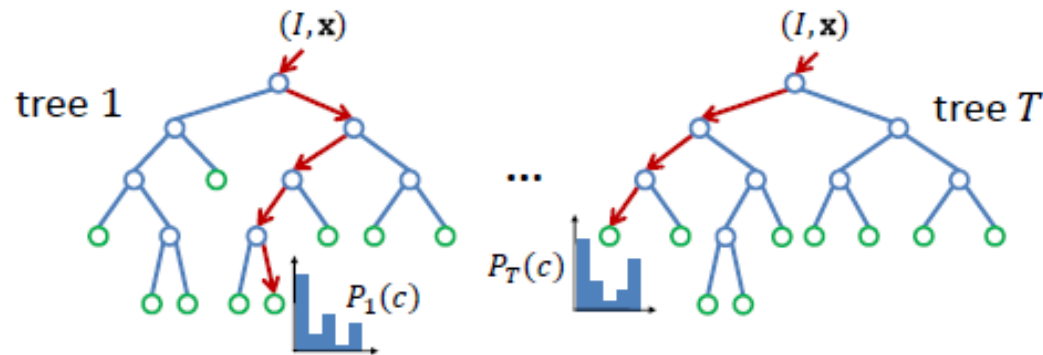depth image ➡ body parts ➡ 3D joint proposals

# Kinect Body Part Recognition

- Features: depth differences between pairs of pixels



- Classification: forest of decision trees

# Large Scale Machine Learning

- Big data

    - Large number of data instances

    - Large number of features


- Solution: distribute computation (parallel computation)

    - GPU (Graphics Processing Unit)

    - Many cores

UNIVERSITY OF
**WATERLOO**

# GPU computation

- Many Machine Learning algorithms consist of vector, matrix and tensor operations

  - A tensor is a multidimensional array

- GPU (Graphics Processing Units) can perform arithmetic operations on all elements of a tensor in parallel

- Packages that facilitate ML programming on GPUs: Keras, PyTorch, TensorFlow, MXNet, Theano, Caffe, DL4J
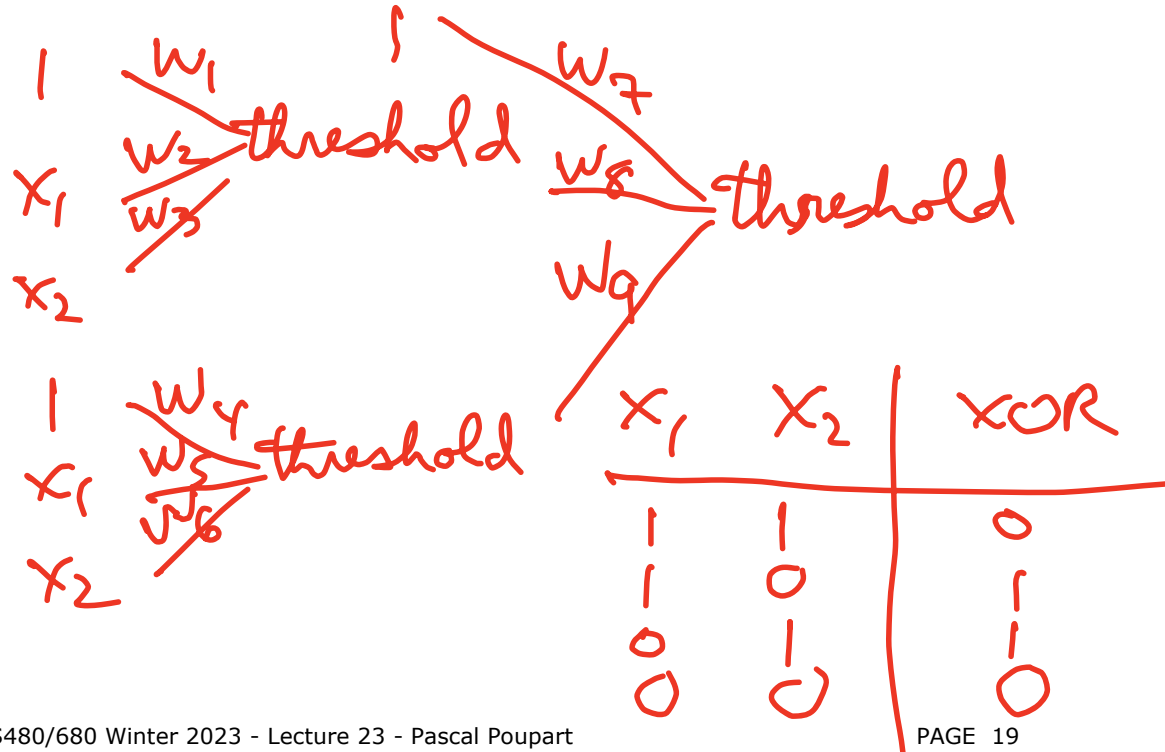
UNIVERSITY OF
WATERLOO

# Multicore Computation

- Idea: Train a different classifier/predictor with a subset of the data on each core


- How can we combine the classifiers/predictors?

- Should we take the average of the parameters of the classifiers/predictors?


**No**, this might lead to a worse classifier/predictor. This is especially problematic for models with hidden variables/units such as neural networks and hidden Markov models

UNIVERSITY OF
**WATERLOO**

# Bad case of parameter averaging

- Consider two threshold neural nets that encode the exclusive-or Boolean function
- Averaging the weights yields new neural net that does not encode exclusive-or



$W_1 = -0.5$ | $-0.5$ | $-0.5$
$W_2 = 1$ | $-1$ | $-$
$W_3 = -1$ | $1$ | $0$
$W_4 = -0.5$ | $-0.5$ | $-0.5$
$W_5 = -1$ | $1$ | $0$
$W_6 = 1$ | $-1$ | $0$
$W_7 = -0.5$ | $-0.5$ | $-0.5$
$W_8 = 1$ | $1$ | $1$
$W_9 = 1$ | $1$ | $1$

| $X_1$ | $X_2$ | XOR |
|---|---|---|
| 1 | 1 | 0 |
| 1 | 0 | 1 |
| 0 | 1 | 1 |
| 0 | 0 | 0 |

UNIVERSITY OF
WATERLOO

# Safely Combining Predictions

- A safe approach to ensemble learning is to **combine the predictions** (not the parameters)

- **Classification:** majority vote of the classes predicted by the classifiers

- **Regression:** average of the predictions computed by the regressors

# Knowledge Distillation

- Technique to train a **small student network $\widetilde{h}$** from a **large teacher network $h$**.

  - Can be used to compress an ensemble of networks into a single network

- Idea: minimize negative log likelihood of target $y$ and cross entropy between teacher and student:

$$\min_{\widetilde{h}} \sum_{(x,y) \in D} \left[ -\log p_{\widetilde{h}}(y|x) - \lambda \sum_{y'} p_h(y'|x) \log p_{\widetilde{h}}(y'|x) \right]$$

UNIVERSITY OF
**WATERLOO**

# Course Perception

- When you have a chance, please fill up the survey and provide feedback about the course (CS480/680) at

https://perceptions.uwaterloo.ca/