# M10. Combined models and Ensemble methods

Manikandan Narayanan

Week 16 (Nov 10- 2025)

PRML Jul-Nov 2025 (Grads Section)

# Acknowledgment of Sources

- Slides based on content from related
  - Courses:
    - IITM – Profs. Arun/Harish/Chandra's PRML offerings (slides, quizzes, notes, etc.), Prof. Ravi's "Intro to ML" slides – cited respectively as [AR], **[HR]**, [CC], [BR] in the bottom right of a slide.
    - India – NPTEL PR course by IISc Prof. PS. Sastry (slides, etc.) – cited as [PSS] in the bottom right of a slide.

  - Books:
    - PRML by Bishop. (content, figures, slides, etc.) – cited as **[CMB]**
    - Pattern Classification by Duda, Hart and Stork. (content, figures, etc.) – [DHS]
    - Mathematics for ML by Deisenroth, Faisal and Ong. (content, figures, etc.) – [DFO]
    - Information Theory, Inference and Learning Algorithms by David JC MacKay – [DJM]

# Outline for Module M10

- M10. Combined models and Ensemble methods
  - **M10.0 Introduction/Motivation**
  - M10.1 Combined models
    - Conditional mixture models
    - Decision trees
  - M10.2 Ensemble methods
    - Parallel ensemble methods (bagging)
    - Sequential ensemble methods (boosting)
  - M10.3 Concluding thoughts

# Machine Learning in Practice

- Real world machine learning problems rarely have a unique and single best solution.

- Multiple thought processes and teams and approaches often yield equally valid but completely different solutions.

- The set of methods for combining many such solutions (classifiers, regressors etc.) into one solution are known as "Ensemble methods"
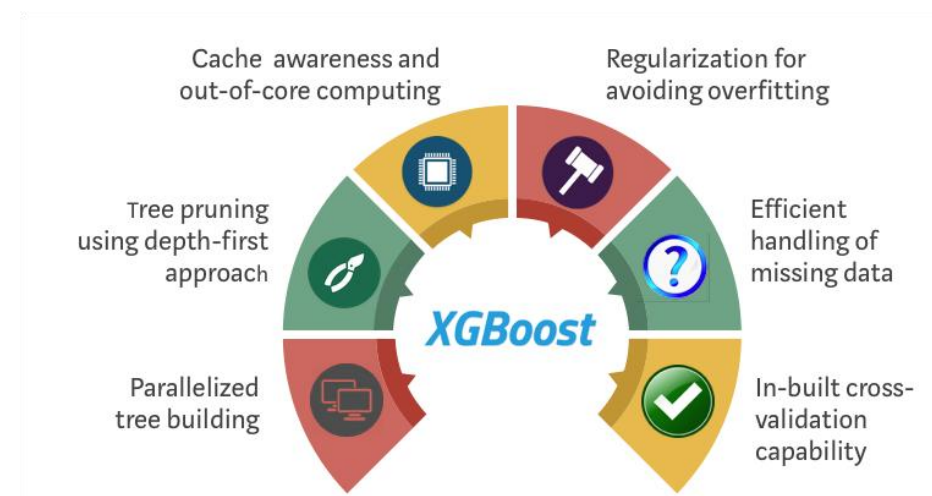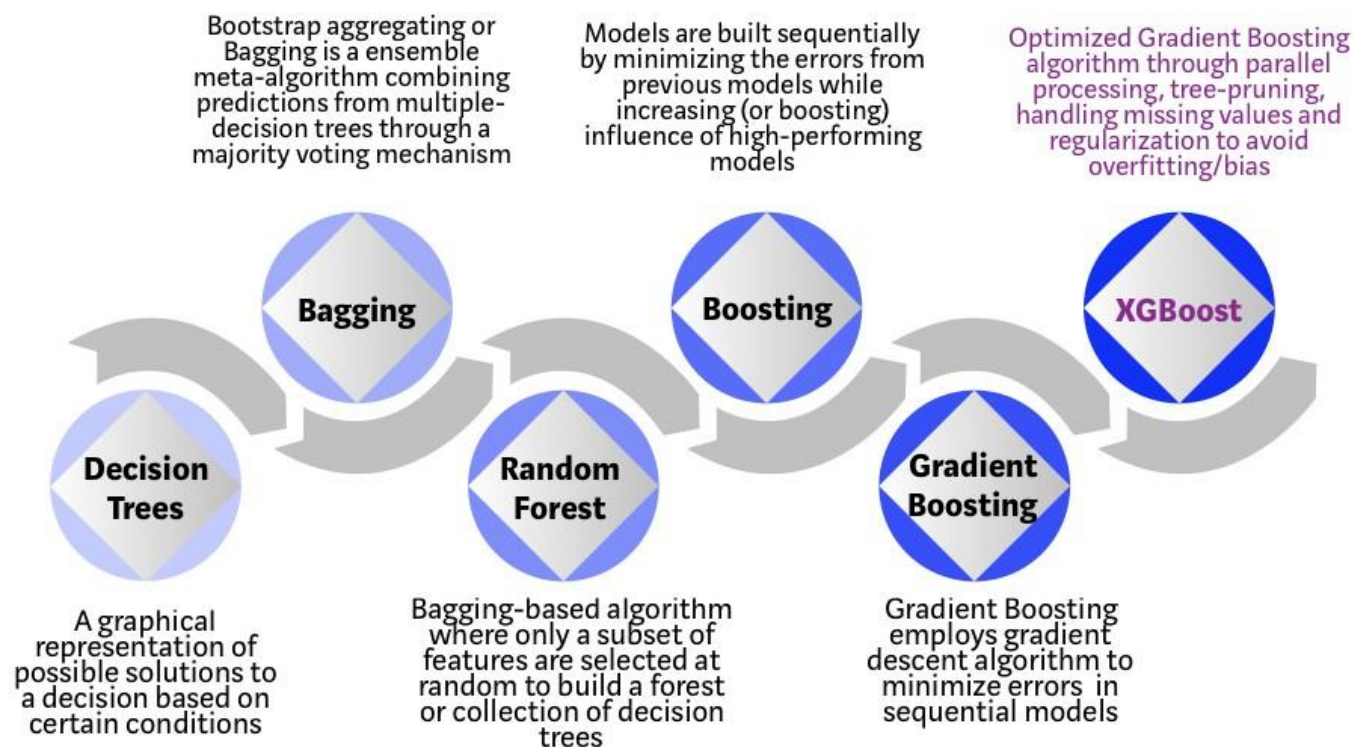
# The Netflix Challenge



- Data set of ~100M star ratings that ~500K users gave to ~18K movies.

- Training data: *<user, movie, date of grade, grade>*

- The grand prize of $1,000,000, to be given to a team which beat Netflix's rating prediction algorithm by 10%

# Netflix challenge: Winning Solution

- The winning team — "BellKor's Pragmatic Chaos" (itself a merger of several teams) combined a total of **107** separate prediction models!!

- The methods used various approaches — factor models, regression models, neighbourhood models, etc.

- Most other teams solutions also included large numbers of disparate models combined together.

- Ensemble methods are almost always the state of the art in any large scale Machine learning problem.
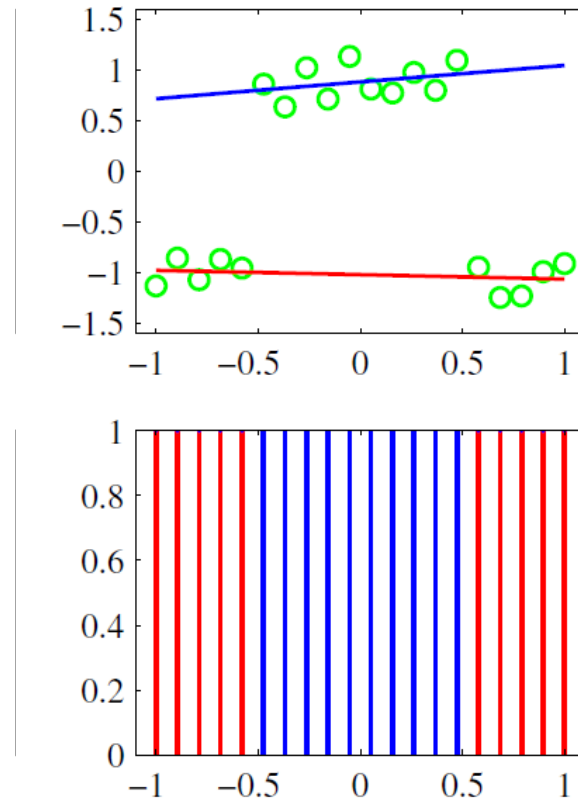
[HR]

# "XGBoost Algorithm: Long May She Reign!"*



Bootstrap aggregating or Bagging is a ensemble meta-algorithm combining predictions from multiple-decision trees through a majority voting mechanism

Models are built sequentially by minimizing the errors from previous models while increasing (or boosting) influence of high-performing models

Optimized Gradient Boosting algorithm through parallel processing, tree-pruning, handling missing values and regularization to avoid overfitting/bias

**Bagging**

**Boosting**

**XGBoost**

**Decision Trees**

**Random Forest**

**Gradient Boosting**

A graphical representation of possible solutions to a decision based on certain conditions

Bagging-based algorithm where only a subset of features are selected at random to build a forest or collection of decision trees

Gradient Boosting employs gradient descent algorithm to minimize errors in sequential models

Cache awareness and out-of-core computing

Regularization for avoiding overfitting

Tree pruning using depth-first approach

Efficient handling of missing data

**XGBoost**

Parallelized tree building

In-built cross-validation capability

# Outline for Module M10

- M10. Combined models and Ensemble methods
    - M10.0 Introduction/Motivation
    - **M10.1 Combined models**
        - **Conditional mixture models**
        - Decision trees
    - M10.2 Ensemble methods
        - Parallel ensemble methods (bagging)
        - Sequential ensemble methods (boosting)
    - M10.3 Concluding thoughts

# Example: Mixtures of linear regression models (mixture of experts)



[CMB]

# Other conditional mixture models

- Mixture of linear classifiers (logistic regression instead of linear regression models): $p(t \mid x) = \sum_k \pi_k p_k(t \mid x)$

- Mixture of experts:
  - Key idea: allow mixing coefficients to be a function of the input x:
  $$p(\mathbf{t}|\mathbf{x}) = \sum_{k=1} \pi_k(\mathbf{x}) p_k(\mathbf{t}|\mathbf{x}).$$

  - Compare with MDN (Mixture Density Network) where all parameters can be input-dependent!
  $$p(\mathbf{t}|\mathbf{x}) = \sum_{k=1} \pi_k(\mathbf{x}) \mathcal{N}\left(\mathbf{t}|\boldsymbol{\mu}_k(\mathbf{x}), \sigma_k^2(\mathbf{x})\right).$$

  - Hierarchical mixture of experts also possible, though we will look at **decision trees** as its hard (non-probab.) version of combining different models:
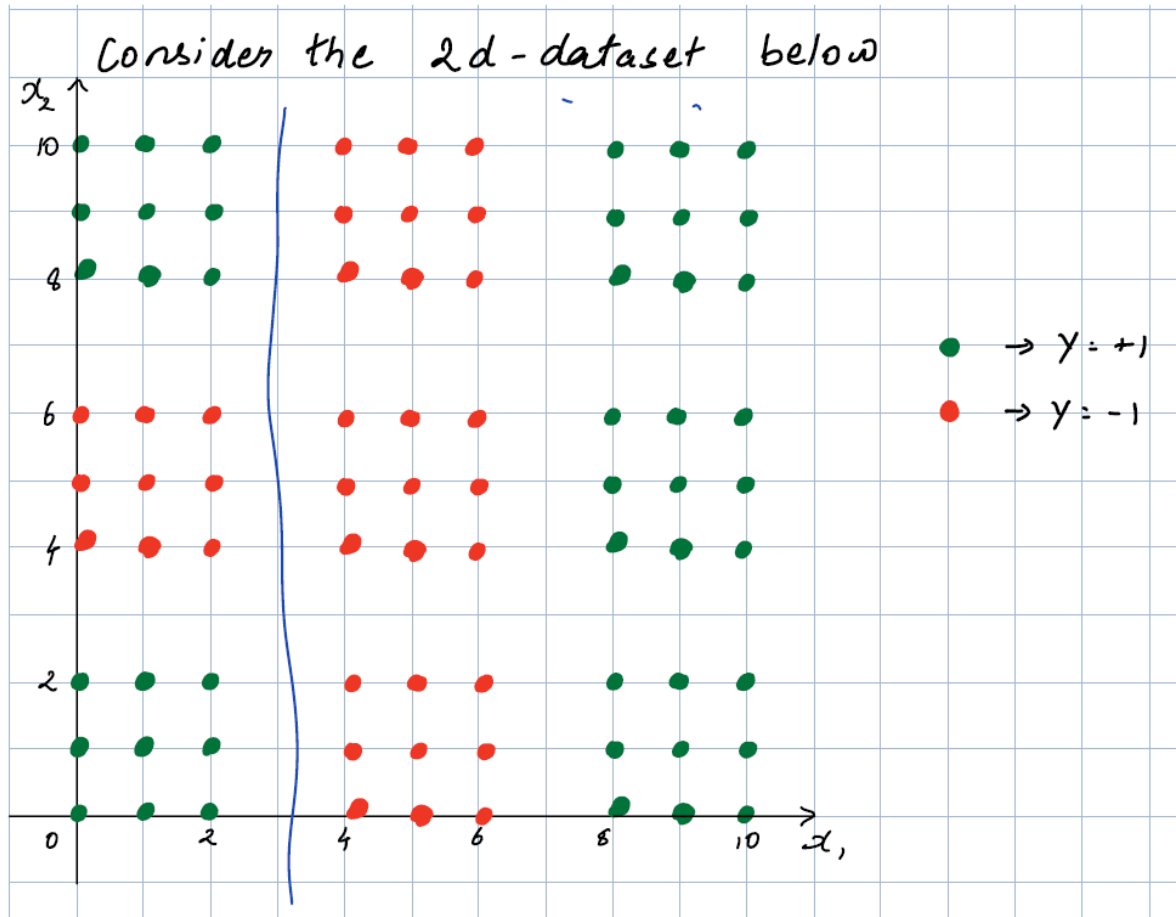
[CMB]

# Outline for Module M10

- M10. Combined models and Ensemble methods
    - M10.0 Introduction/Motivation
    - **M10.1 Combined models**
        - Conditional mixture models
        - **Decision trees**
    - M10.2 Ensemble methods
        - Parallel ensemble methods (bagging)
        - Sequential ensemble methods (boosting)
    - M10.3 Concluding thoughts

# A visual understanding of decision trees

- http://www.r2d3.us/visual-intro-to-machine-learning-part-1/

# Learning a decision tree – toy example



Consider the 2d-dataset below

● → Y = +1
● → Y = -1

1.) What should the Root node be?
   · Evaluate all classifiers of the form on "Entire" Training.

$$h(x) = \begin{cases} +1 & \text{if } x_1 \geq a \\ -1 & \text{if } x_1 < a \end{cases} \quad \text{and}$$

$$h(x) = \begin{cases} +1 & \text{if } x_2 \geq a \\ -1 & \text{if } x_2 < a \end{cases}$$

and their negations.

# Toy example – evaluating all "feature X threshold" combinations at the root node!

## *Accuracy*

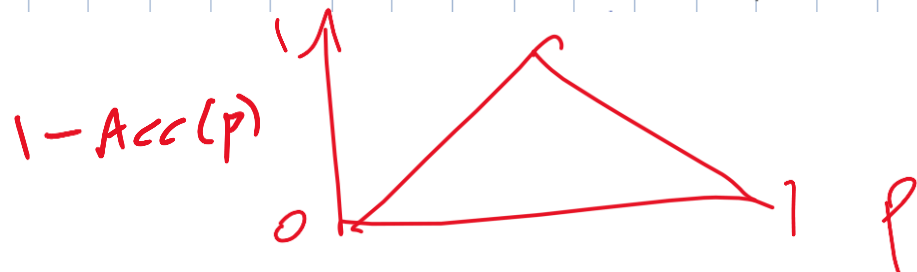For simplicity we will evaluate only 4 such classifiers:

(a) 
$$h(x) = +1 \quad \text{if} \quad x_1 \geq 7$$
$$= -1 \quad \text{if} \quad x_1 < 7$$

$$\text{Accuracy} = \frac{3+4}{9} = \frac{7}{9}$$

(a') 
$$h(x) = -1 \quad \text{if} \quad x_1 \geq 7$$
$$= +1 \quad \quad x_1 \leq 7$$

Let $L: x_1 < 7$
$R: x_1 \geq 7$

$$\text{Accuracy} = \frac{0+2}{9} = \frac{2}{9}$$

$1 - Acc(p)$

## *Entropy*

Define $H(p) = p \log \frac{1}{p} + (1-p) \log \left(\frac{1}{1-p}\right)$  (logarithm base is 2)

Avg. Entropy of split (a) $= P_L H_L + P_R H_R$.

$P_L$ = Fraction of points on the left. e.g. above $P_L = \frac{54}{81}$

$H_L = H(q_L)$

where $q_L$ = Fraction of positive points on the left. e.g above $q_L = \frac{18}{54}$

∴ Entropy of split (a) above

$$= \frac{54}{81} H\left(\frac{18}{54}\right) + \frac{27}{81} H(1)$$

$$= \frac{2}{3} H\left(\frac{1}{3}\right) + \frac{1}{3} H(1) = \frac{2}{3} H\left(\frac{1}{3}\right)$$

$$= \frac{2}{3} \left(\frac{1}{3} \log 3 + \frac{2}{3} \log \frac{3}{2}\right) = \frac{2}{3} \left(\log 3 - \frac{2}{3}\right)$$

$H(p)$

[HR]

# Exercise: try for 3 other splits (& their negations) to find the best root node split.

- E.g.,

(b)    $h(x) = +1$    if    $x_1 \geq 3$
       $= -1$    if    $x_1 < 3$

Accuracy $= \dfrac{3+1}{9} = \dfrac{4}{9}$

lllrly accuracy of (b') which is the negation of above

$= 1 - \dfrac{4}{9} = \dfrac{5}{9}$

$L: x_1 < 3$
$R: x_1 \geq 3$

Entropy of split b : $P_L H_L + P_R H_R$

$P_L = \dfrac{1}{3}$    $q_L = \dfrac{2}{3}$    : Entropy : $\dfrac{1}{3} H\left(\dfrac{2}{3}\right) + \dfrac{2}{3} H\left(\dfrac{1}{2}\right)$

$P_R = \dfrac{2}{3}$    $q_R = \dfrac{1}{2}$

Clearly   Option (a)   is   best for root node:

(In this case using either max accuracy or min entropy gives same answer. But not always)

$x_1 \geq 7$

No    Yes

Recurse    Stop

[HR]

# Exercise: Recursively solve subproblems on left and right to derive this final tree
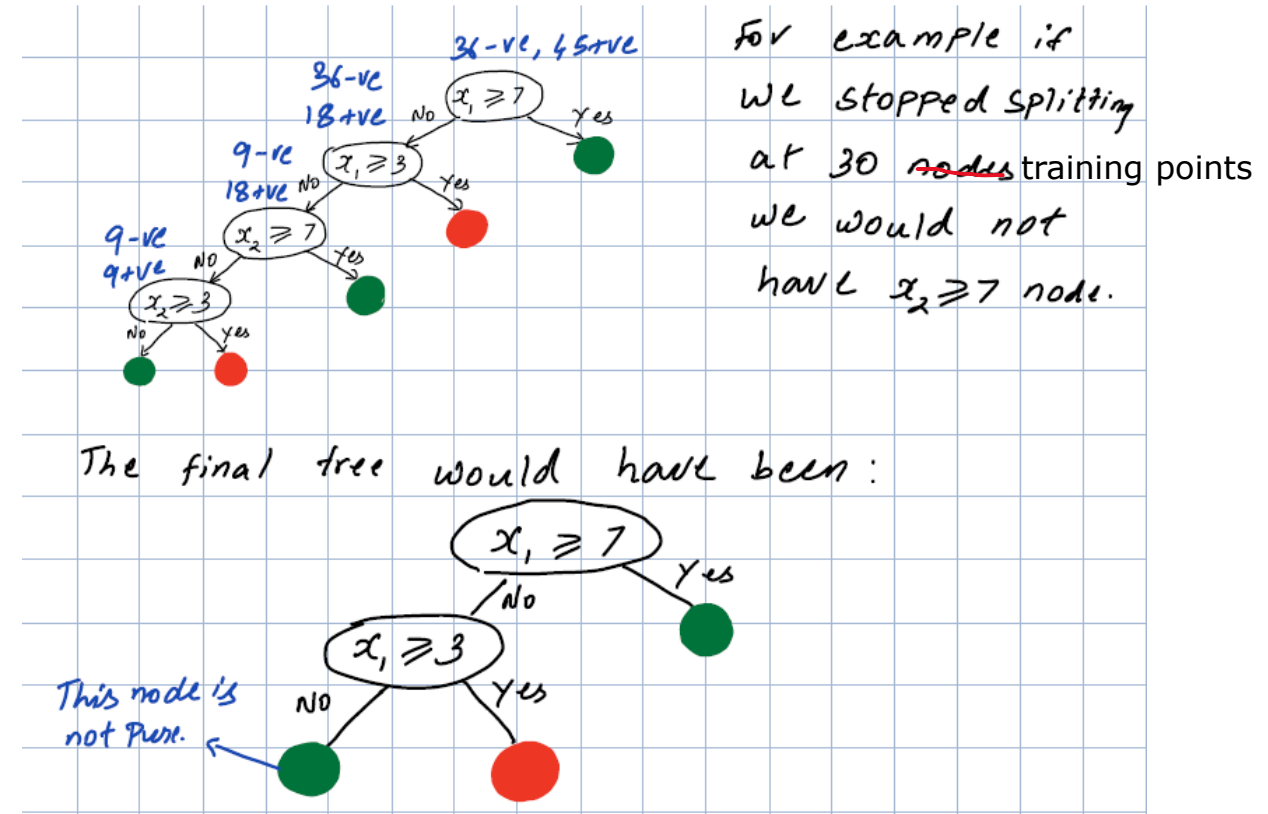


(Pure node)

final tree

# Regularization of a decision tree

- Goal: Bring down the # of nodes in the decision tree without losing too much on accuracy.

- Solution: Stop early...
    i) ...at a certain depth of the tree, or
    ii) ...when number of training points is less than some number.

*preferred*

# Outline for Module M10

- M10. Combined models and Ensemble methods
    - M10.0 Introduction/Motivation
    - M10.1 Combined models
        - Conditional mixture models
        - Decision trees
    - **M10.2 Ensemble methods or Committee models**
        - **Parallel ensemble methods (bagging)**
        - **Sequential ensemble methods (boosting)**
    - M10.3 Concluding thoughts

# Committees

- Committees are ensemble methods that average the predictions of many individual learners

- Two very different approaches:
  - Bagging – average of **parallel**y/separately-trained **high-capacity** learners
  - Boosting – average of **sequential**ly/adaptively-trained **weak** learners

- Bias-variance decomposition helps in understanding certain aspects of bagging, and computational/statistical learning theory helps derive certain performance bounds on boosting.
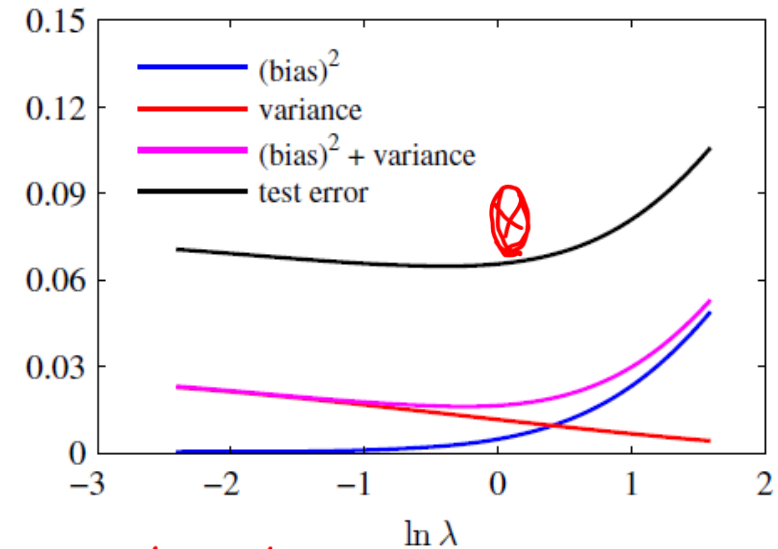
# Recall: Bias-variance analysis summary

expected loss $= (\text{bias})^2 + \text{variance} + \text{noise}$

where

$$(\text{bias})^2 = \int \{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2 p(\mathbf{x}) \, d\mathbf{x}$$

$$\text{variance} = \int \mathbb{E}_{\mathcal{D}}\left[\{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}^2\right] p(\mathbf{x}) \, d\mathbf{x}$$

$$\text{noise} = \int \{h(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) \, d\mathbf{x} \, dt$$
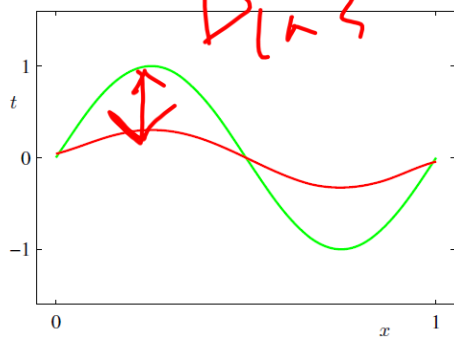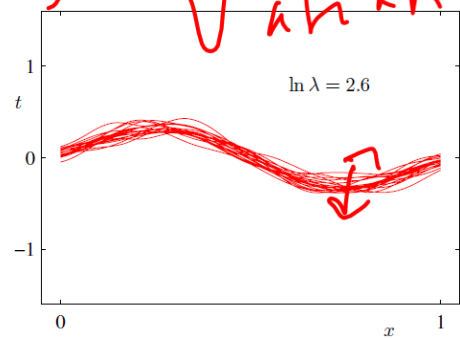
# Recall: Bias-variance in pictures



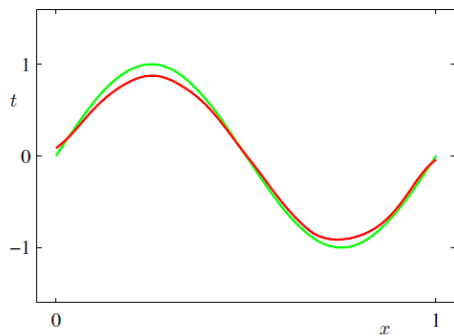cplxty

lo

cnt.

hi

Variance

Bias

$\ln \lambda = 2.6$

$\ln \lambda = -0.31$

$\ln \lambda = -2.4$

[CMB]
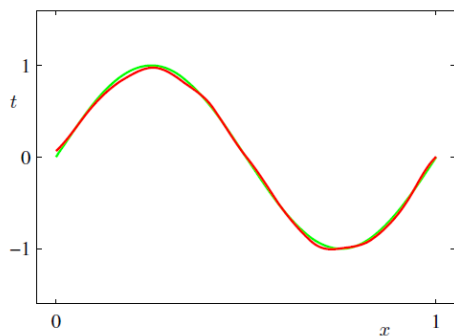
# Switch to Harish's slides on Ensemble methods

[ https://drive.google.com/file/d/1KUy2mziZ1pz-A0NeqlAV_qTP5Rq47Kkt/view?usp=sharing ]

See annotated version of above called "N11_ensemble_methods.ExtraNotes.pdf" in moodle.

# Outline for Module M10

- M10. Combined models and Ensemble methods
  - M10.0 Introduction/Motivation
  - M10.1 Combined models
    - Conditional mixture models
    - Decision trees
  - M10.2 Ensemble methods
    - Parallel ensemble methods (bagging)
    - Sequential ensemble methods (boosting)
  - **M10.3 Concluding thoughts**

# "XGBoost Algorithm: Long May She Reign!"*



Bootstrap aggregating or Bagging is a ensemble meta-algorithm combining predictions from multiple-decision trees through a majority voting mechanism

Models are built sequentially by minimizing the errors from previous models while increasing (or boosting) influence of high-performing models

Optimized Gradient Boosting algorithm through parallel processing, tree-pruning, handling missing values and regularization to avoid overfitting/bias

**Bagging**

**Boosting**

**XGBoost**

**Decision Trees**

A graphical representation of possible solutions to a decision based on certain conditions

**Random Forest**

Bagging-based algorithm where only a subset of features are selected at random to build a forest or collection of decision trees

**Gradient Boosting**

Gradient Boosting employs gradient descent algorithm to minimize errors in sequential models

Cache awareness and out-of-core computing

Regularization for avoiding overfitting

Tree pruning using depth-first approach

Efficient handling of missing data

**XGBoost**

Parallelized tree building

In-built cross-validation capability

# Concluding thoughts & next steps

**Supervised learning:** Fun starts when you can take the different classification or regression models you've learned from this course, and apply your knowledge to choose the right model or right method of combining models in a systematic rather than brute-force fashion!

- Unified view of different models helps towards above goal – fixed vs. selective (SVM) vs. adaptive (ANN) basis functions; loss fn. view of different classifiers, etc.
- Understanding conceptual/mathematical foundations of different methods – beyond popular "blog" descriptions -- also helps towards above goal.

similar also applies for **Unsupervised learning:** spectral clustering and dimensionality reduction (PCA) both involved understanding the eigenspectrum of a matrix (optimizing $u^T A u$).

# Thank you!