

M10. Combined models and Ensemble methods

B Ravindran

PRML Jul-Aug 2021 (BR 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
 - M10.3 Ensemble methods
 - Parallel ensemble methods (bagging)
 - Sequential ensemble methods (boosting)
 - M10.4 Parting 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

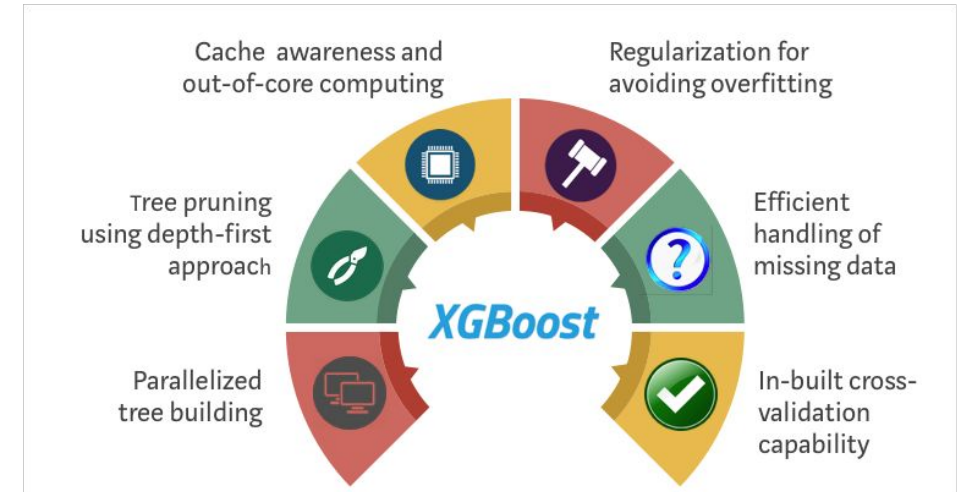
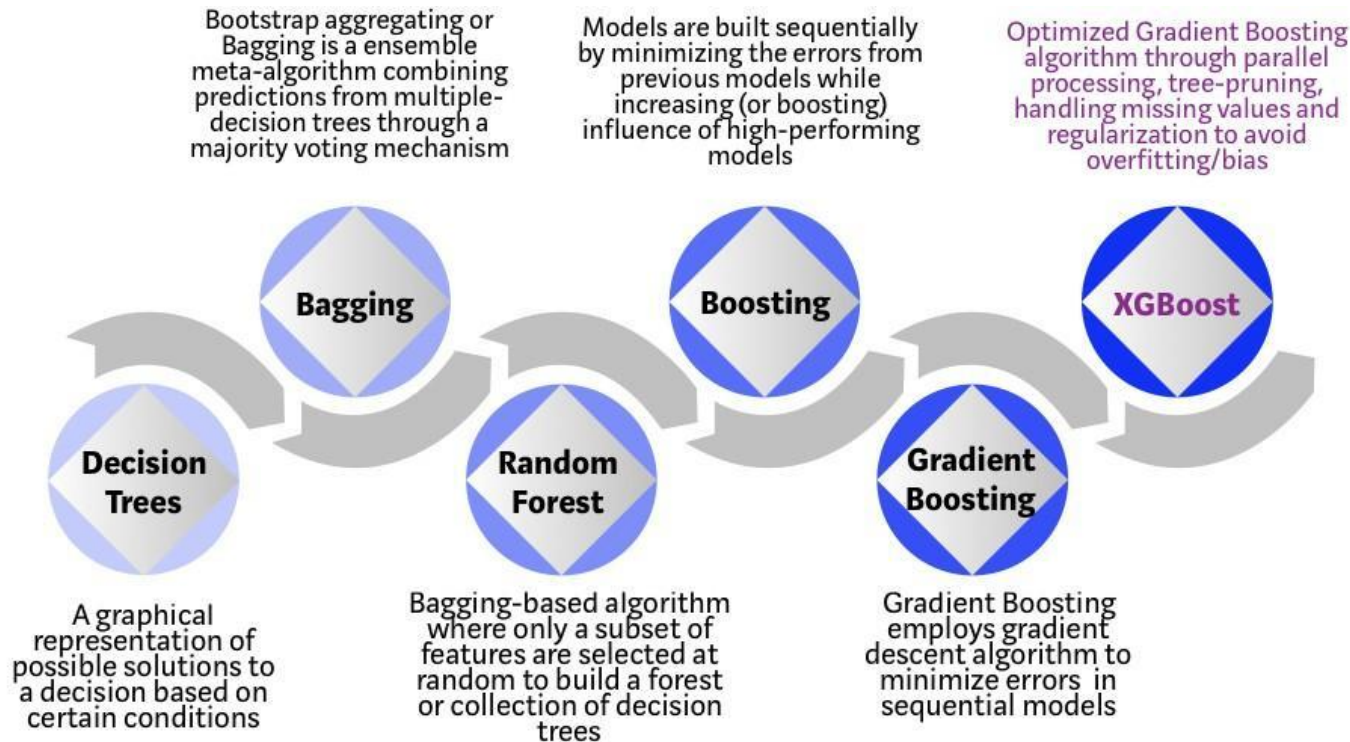


- Data set of ~100M star ratings that ~500K users gave to ~18K movies.
- Training data: $\langle user, movie, date\ of\ grade, grade \rangle$
- The grand prize of \$1,000,000, to be given to a team which beat Netflix's rating prediction algorithm by 10%

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.

“XGBoost Algorithm: Long May She Reign!”*

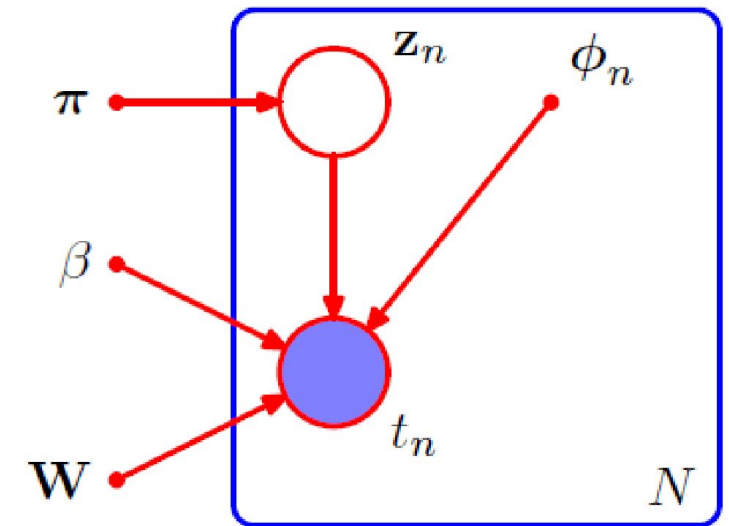


Outline for Module M10

- M10. Combined models and Ensemble methods
 - M10.0 Introduction/Motivation
 - **M10.1 Combined models**
 - **Conditional mixture models**
 - Decision trees
 - M10.3 Ensemble methods
 - Parallel ensemble methods (bagging)
 - Sequential ensemble methods (boosting)
 - M10.4 Parting thoughts

Conditional mixture models: mixture of linear regression models

$$p(t|\boldsymbol{\theta}) = \sum_{k=1}^K \pi_k \mathcal{N}(t | \mathbf{w}_k^T \boldsymbol{\phi}, \beta^{-1})$$



$$\gamma_{nk} = \mathbb{E}[z_{nk}] = p(k | \phi_n, \boldsymbol{\theta}^{\text{old}}) = \frac{\pi_k \mathcal{N}(t_n | \mathbf{w}_k^T \boldsymbol{\phi}_n, \beta^{-1})}{\sum_j \pi_j \mathcal{N}(t_n | \mathbf{w}_j^T \boldsymbol{\phi}_n, \beta^{-1})}.$$

MLE using EM algo.

- Maximize (Marginal) Likelihood: $\ln p(\mathbf{t}|\boldsymbol{\theta}) = \sum_{n=1}^N \ln \left(\sum_{k=1}^K \pi_k \mathcal{N}(t_n | \mathbf{w}_k^T \boldsymbol{\phi}_n, \beta^{-1}) \right)$

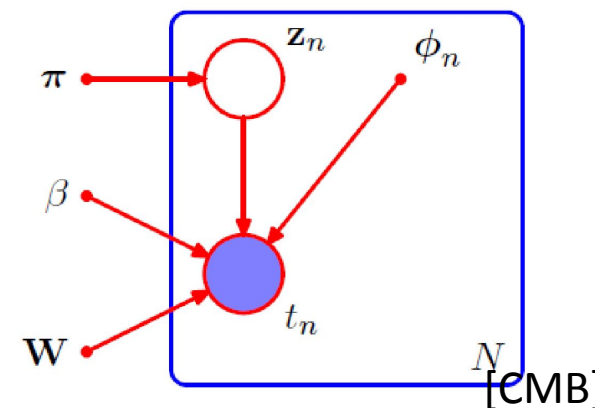
- E-step: Compute responsibilities γ_{nk} .

- M-step: Maximize $Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{old}}) = \mathbb{E}_{\mathbf{Z}} [\ln p(\mathbf{t}, \mathbf{Z}|\boldsymbol{\theta})] = \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} \{ \ln \pi_k + \ln \mathcal{N}(t_n | \mathbf{w}_k^T \boldsymbol{\phi}_n, \beta^{-1}) \}$.

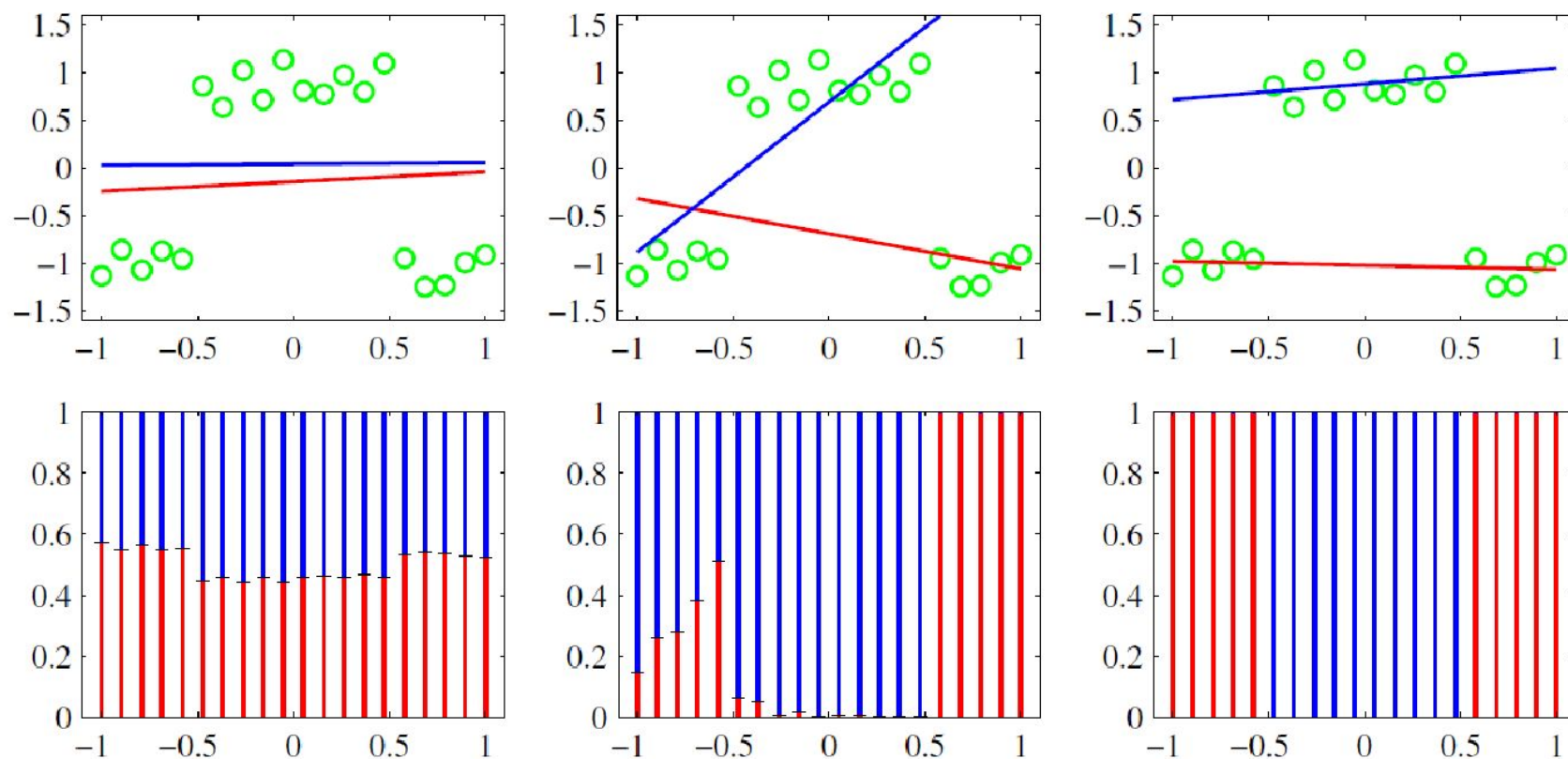
$$\pi_k = \frac{1}{N} \sum_{n=1}^N \gamma_{nk}.$$

$$\mathbf{w}_k = (\boldsymbol{\Phi}^T \mathbf{R}_k \boldsymbol{\Phi})^{-1} \boldsymbol{\Phi}^T \mathbf{R}_k \mathbf{t}. \quad \mathbf{R}_k = \text{diag}(\gamma_{nk})$$

$$\frac{1}{\beta} = \frac{1}{N} \sum_{n=1}^N \sum_{k=1}^K \gamma_{nk} (t_n - \mathbf{w}_k^T \boldsymbol{\phi}_n)^2.$$



Example: Mixtures of linear regression models



Other conditional mixture models

- Mixture of linear classifiers (logistic regression instead of linear regression models)

- Mixture of experts:

- Key idea: allow mixing coefficients to be a function of the input \mathbf{x} :

$$p(\mathbf{t}|\mathbf{x}) = \sum_{k=1}^K \pi_k(\mathbf{x}) p_k(\mathbf{t}|\mathbf{x}).$$

- Compare with MDNs where all parameters can be input-dependent!

$$p(\mathbf{t}|\mathbf{x}) = \sum_{k=1}^K \pi_k(\mathbf{x}) \mathcal{N}(\mathbf{t} | \boldsymbol{\mu}_k(\mathbf{x}), \sigma_k^2(\mathbf{x})).$$

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

Outline for Module M10

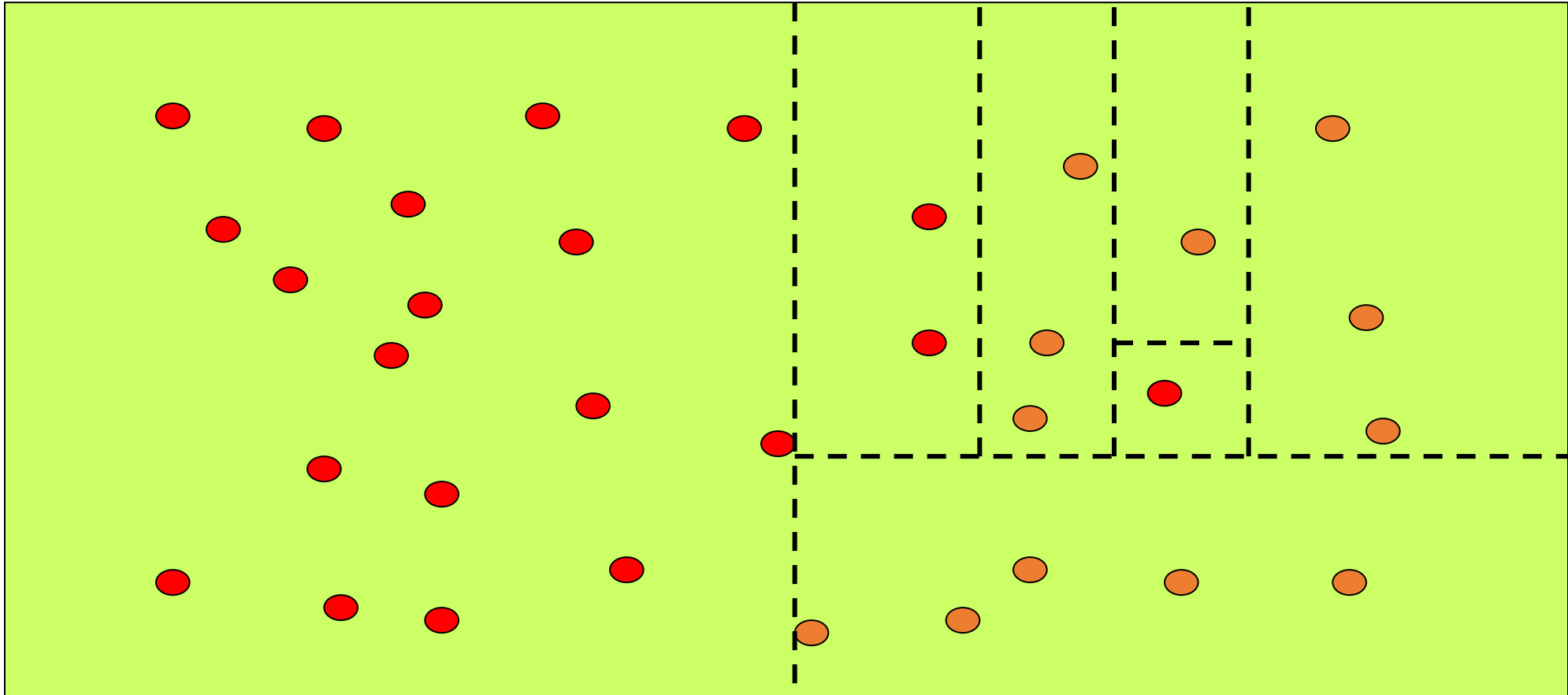
- M10. Combined models and Ensemble methods
 - M10.0 Introduction/Motivation
 - **M10.1 Combined models**
 - Conditional mixture models
 - **Decision trees**
 - M10.3 Ensemble methods
 - Parallel ensemble methods (bagging)
 - Sequential ensemble methods (boosting)
 - M10.4 Parting thoughts

A visual understanding of decision trees

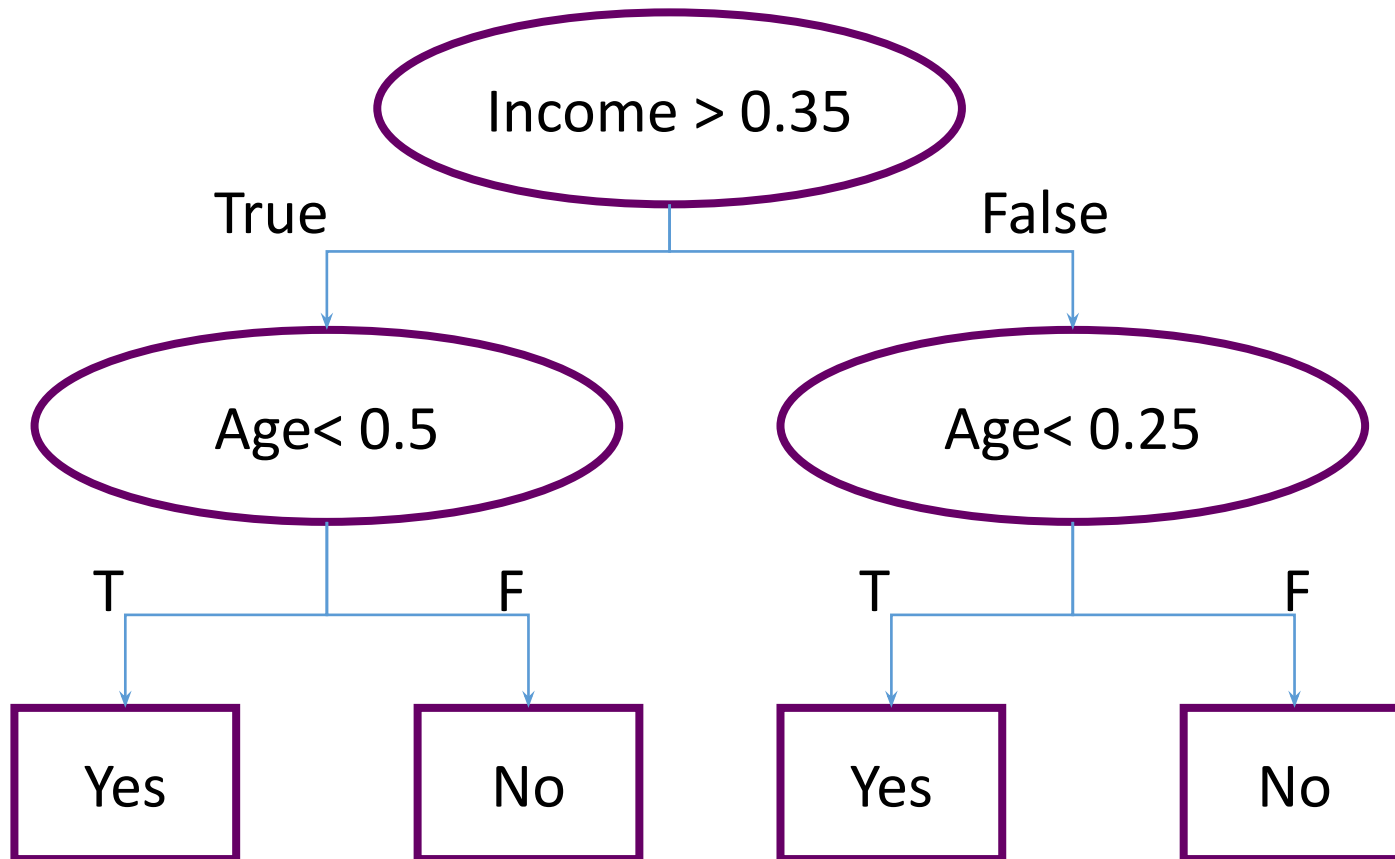
- <http://www.r2d3.us/visual-intro-to-machine-learning-part-1/>

Decision Trees

Decision boundaries



A Decision Tree



If $\text{Income} > 0.35$ and
 $\text{Age} < 0.5$
then $\text{Buys} = \text{Yes}$

If $\text{Income} > 0.35$ and
 $\text{Age} \geq 0.5$
then $\text{Buys} = \text{No}$

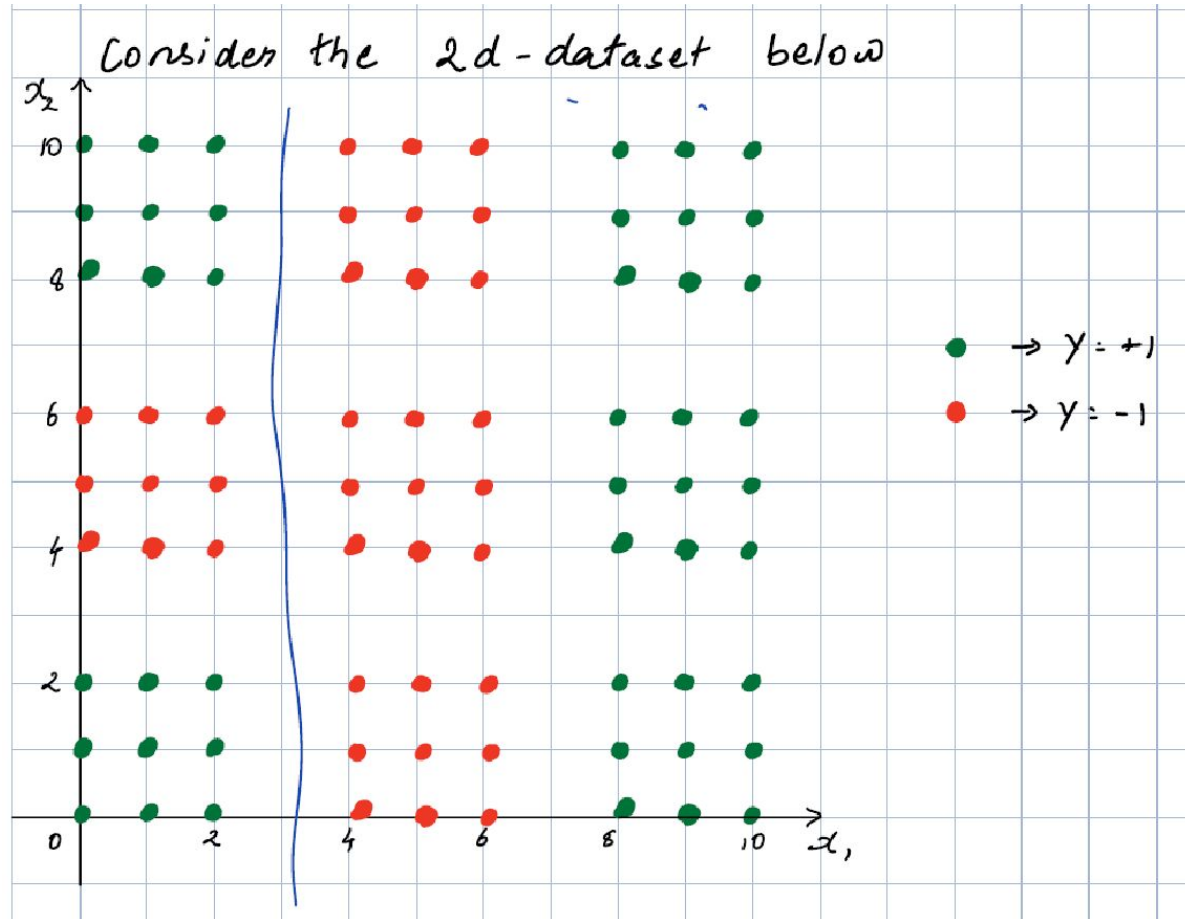
If $\text{Income} \leq 0.35$ and
 $\text{Age} < 0.25$
then $\text{Buys} = \text{Yes}$

If $\text{Income} \leq 0.35$ and
 $\text{Age} \geq 0.25$
then $\text{Buys} = \text{No}$

Constructing decision trees

- Strategy: top down
Recursive *divide-and-conquer* fashion
 - First: select attribute for root node
Create branch for each possible attribute value
 - Then: split instances into subsets
One for each branch extending from the node
 - Finally: repeat recursively for each branch, using only instances that reach the branch
- Stop if all instances have the same class

Learning a decision tree – toy example



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:

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

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

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

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

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

Entropy

Define $H(P) = P \log \frac{1}{P} + (1-P) \log \frac{1}{1-P}$ (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 $P_L = \frac{54}{81}$
 $H_L = H(P_L)$
Where q_L = Fraction of positive points on the left. e.g. above $q_L = \frac{18}{54}$

$$\begin{aligned} \therefore \text{Entropy of split (a) above} &= \frac{54}{81} H\left(\frac{18}{54}\right) + \frac{27}{81} H(1) \\ &= \frac{2}{3} H\left(\frac{1}{3}\right) + \frac{1}{3} H(1) = \frac{2}{3} H\left(\frac{1}{3}\right) \\ &= \frac{2}{3} \left(\frac{1}{3} \log 3 + \frac{2}{3} \log \frac{3}{2} \right) = \frac{2}{3} \left(\log 3 - \frac{2}{3} \right) \end{aligned}$$

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

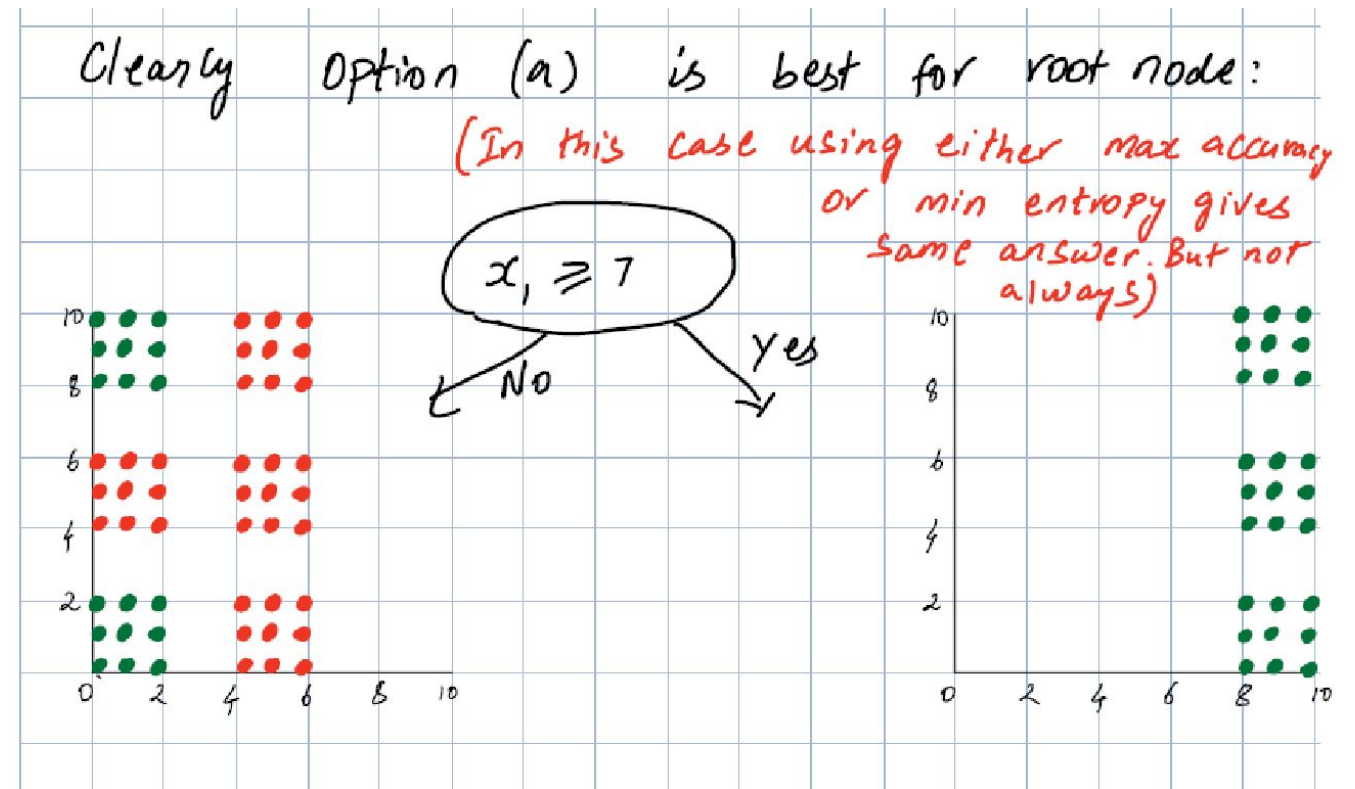
• E.g.,

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

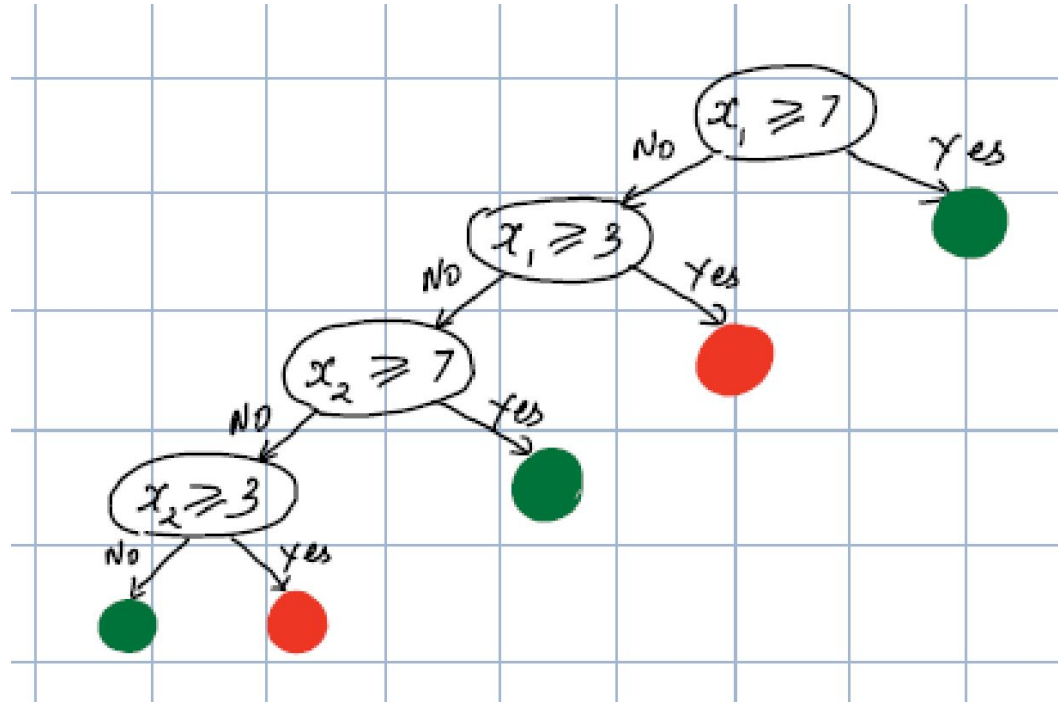
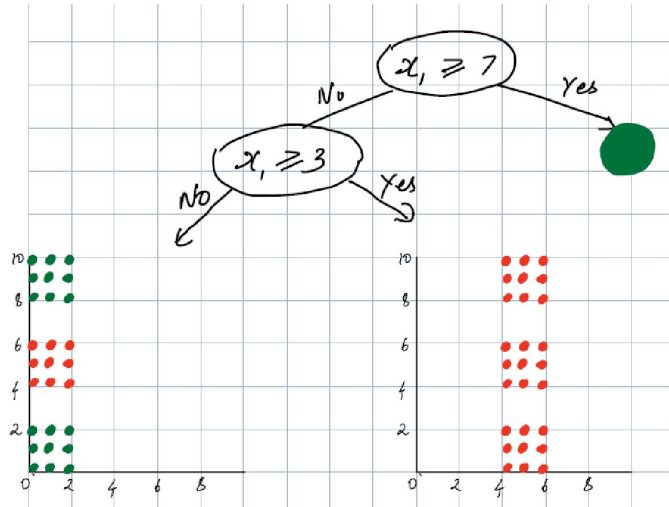
Accuracy = $\frac{3+1}{9} = \frac{4}{9}$

Only accuracy of (b') which is the negation of above
 $= 1 - \frac{4}{9} = \frac{5}{9}$

Entropy of split b : $P_L H_L + P_R H_R$
 $P_L = \frac{1}{3}$ $q_L = \frac{2}{3}$: Entropy : $\frac{1}{3} H\left(\frac{2}{3}\right) + \frac{2}{3} H\left(\frac{1}{3}\right)$
 $P_R = \frac{2}{3}$ $q_R = \frac{1}{3}$ L: $x_1 < 3$ R: $x_1 \geq 3$



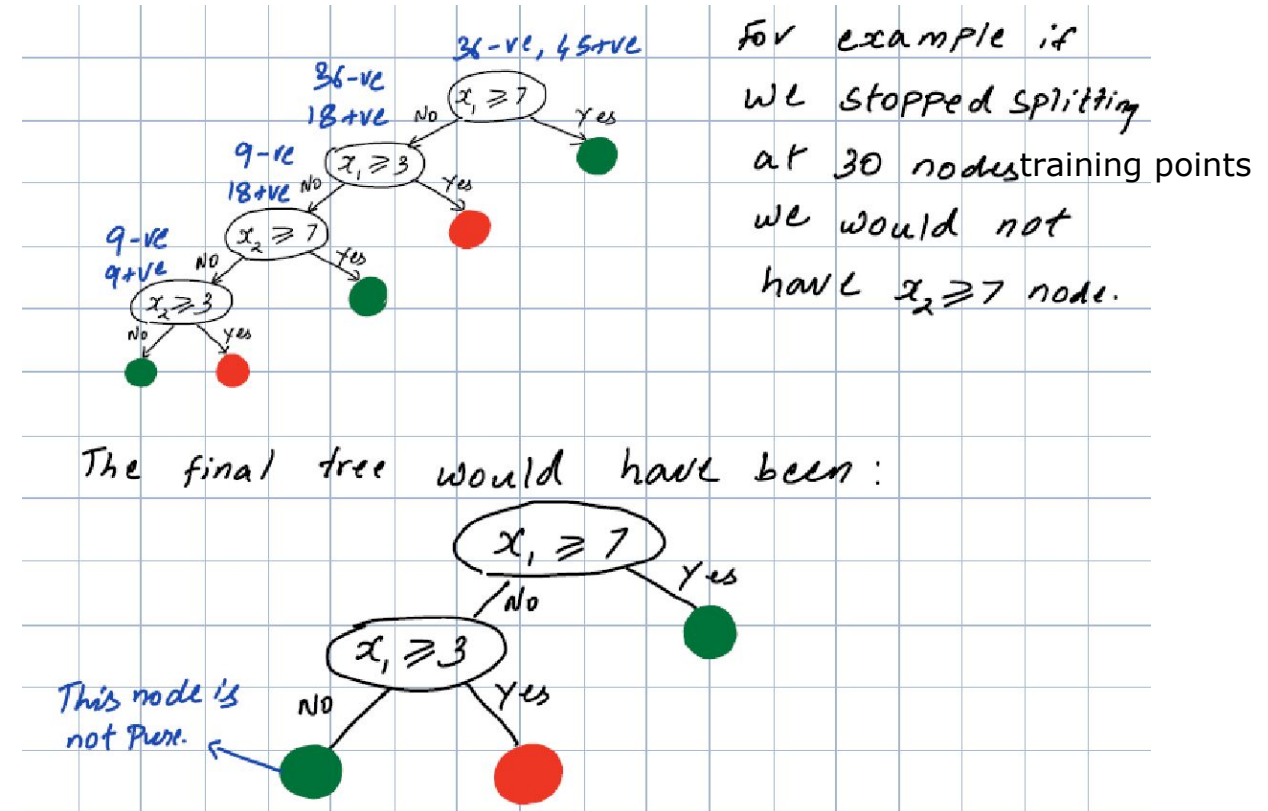
Exercise: Recursively solve subproblems on left and right to derive this final tree



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.



Outline for Module M10

- M10. Combined models and Ensemble methods
 - M10.0 Introduction/Motivation
 - M10.1 Combined models
 - Conditional mixture models
 - Decision trees
 - **M10.3 Ensemble methods or Committee models**
 - Parallel ensemble methods (bagging)
 - Sequential ensemble methods (boosting)
 - M10.4 Parting thoughts

Committees

- Committees are ensemble methods that average the predictions of many individual learners
- Two very different approaches:
 - Bagging – average of **parallelly**/separately-trained **high-capacity** learners
 - Boosting – average of **sequentially**/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.

Bias-variance analysis

$$h(\mathbf{x}) = \mathbb{E}[t|\mathbf{x}] = \int tp(t|\mathbf{x}) dt.$$

$$\mathbb{E}[L] = \int \{y(\mathbf{x}) - h(\mathbf{x})\}^2 p(\mathbf{x}) d\mathbf{x} + \int \{h(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) d\mathbf{x} dt.$$

$$\{y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x})\}^2.$$

$$\begin{aligned} & \{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] + \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2 \\ &= \{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}^2 + \{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2 \\ & \quad + 2\{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}\{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}. \end{aligned}$$

$$\begin{aligned} & \mathbb{E}_{\mathcal{D}} [\{y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x})\}^2] \\ &= \underbrace{\{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2}_{(\text{bias})^2} + \underbrace{\mathbb{E}_{\mathcal{D}} [\{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}^2]}_{\text{variance}}. \end{aligned}$$

Bias-variance analysis summary

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

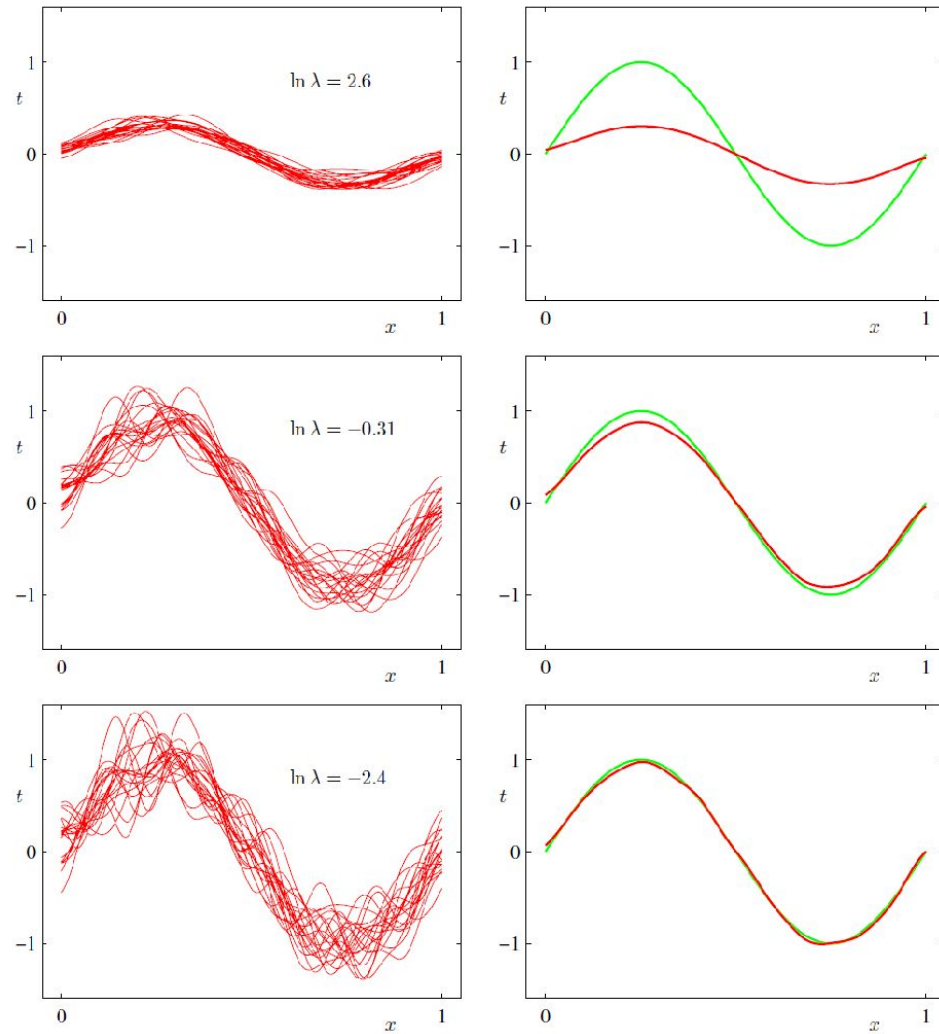
where

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

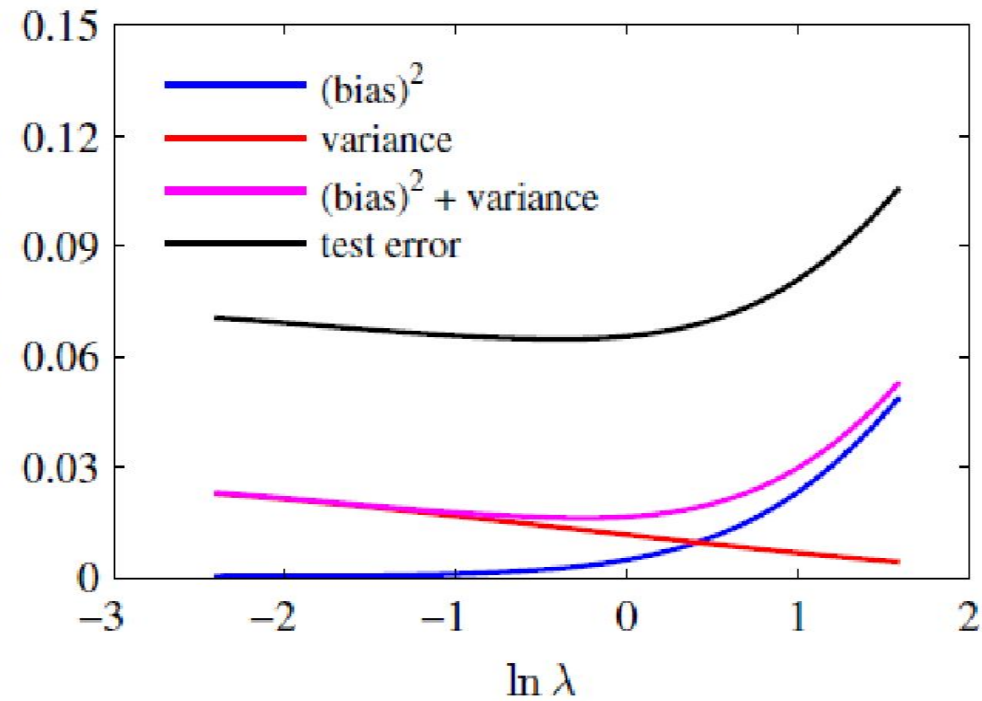
$$\text{variance} = \int \mathbb{E}_{\mathcal{D}} [\{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}^2] p(\mathbf{x}) \, d\mathbf{x}$$

$$\text{noise} = \int \{h(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) \, d\mathbf{x} \, dt$$

Bias-variance in pictures



Bias-variance analysis: empirical view



Bias-variance analysis: applicability in practice?

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

where

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

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

$$\text{noise} = \int \{h(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) \, d\mathbf{x} \, dt$$

$$\bar{y}(x) = \frac{1}{L} \sum_{l=1}^L y^{(l)}(x)$$

$$(\text{bias})^2 = \frac{1}{N} \sum_{n=1}^N \{\bar{y}(x_n) - h(x_n)\}^2$$

$$\text{variance} = \frac{1}{N} \sum_{n=1}^N \frac{1}{L} \sum_{l=1}^L \{y^{(l)}(x_n) - \bar{y}(x_n)\}^2$$

Recall: Bias-variance analysis summary

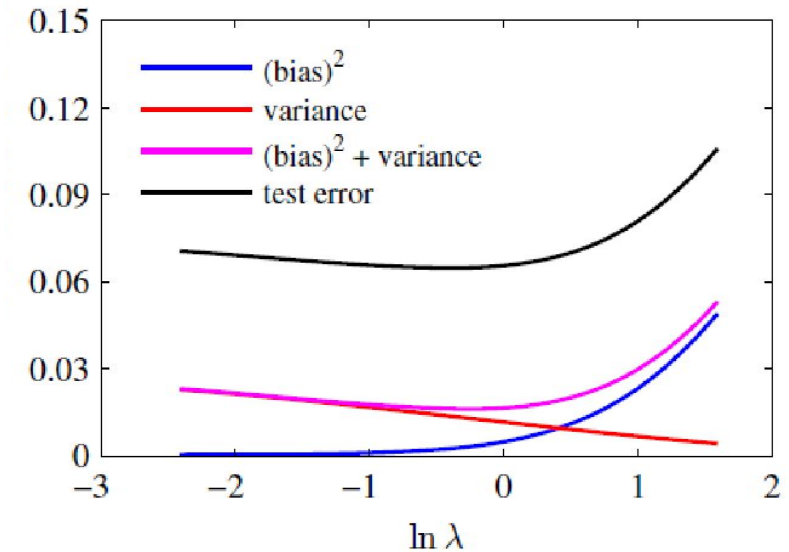
$$\text{expected loss} = (\text{bias})^2 + \text{variance} + \text{noise}$$

where

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

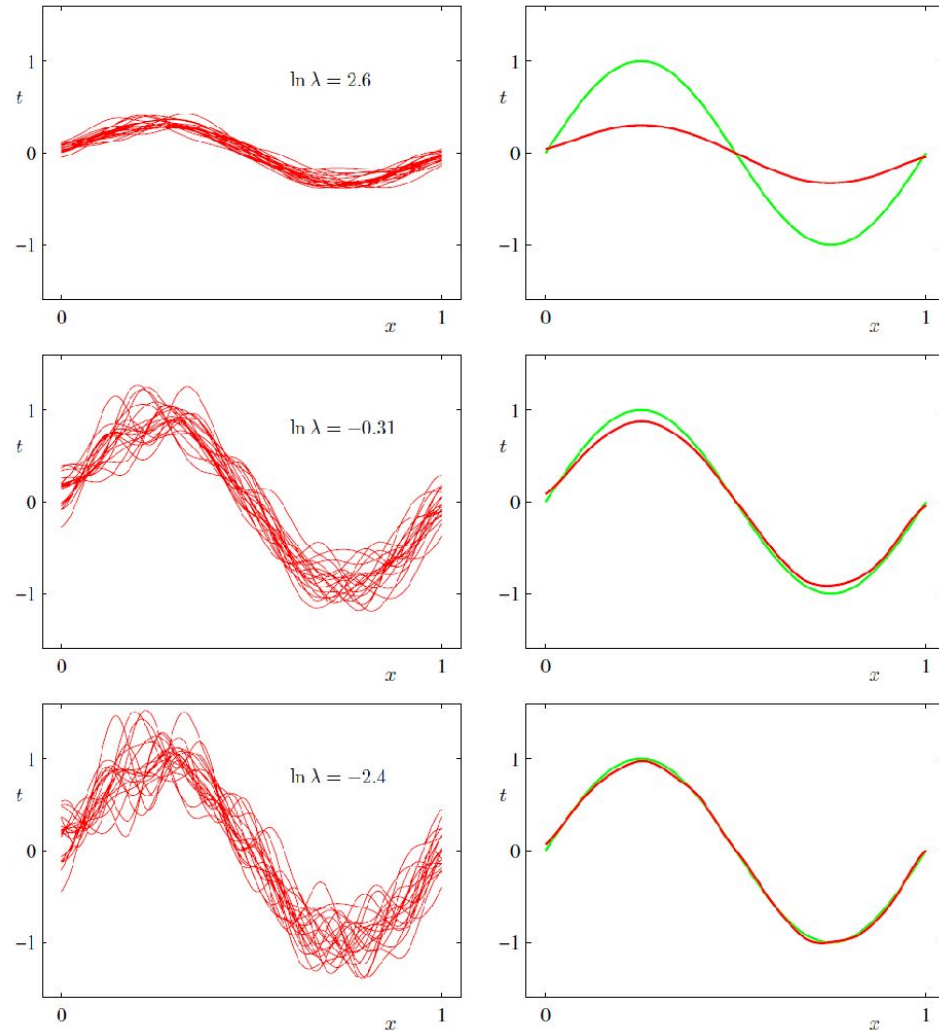
$$\text{variance} = \int \mathbb{E}_{\mathcal{D}} [\{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}^2] p(\mathbf{x}) d\mathbf{x}$$

$$\text{noise} = \int \{h(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) d\mathbf{x} dt$$



Exercise: cf. slides below or zoom quiz for careful understanding of what random variables the expectation above is taken over!

Recall: Bias-variance in pictures



Switch to Harish's slides on Ensemble methods

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

Outline for Module M10

- M10. Combined models and Ensemble methods
 - M10.0 Introduction/Motivation
 - M10.1 Combined models
 - Conditional mixture models
 - Decision trees
 - M10.3 Ensemble methods
 - Parallel ensemble methods (bagging)
 - Sequential ensemble methods (boosting)
 - **M10.4 Parting thoughts**

Recall: There and back again: Planned syllabus (tentative)

Subset of topics below to be covered (not necessarily in the same order):

1. Overview of PR/ML problems/algorithms (PR tasks/systems, ML paradigms)
2. Bayesian decision theory (Bayes classifier, loss functions)
3. Density estimation (Maximum likelihood, Bayesian estimation, Expectation Maximization (EM) for mixture density estimation, Non-parametric methods)
4. Linear models for classification and regression (Linear discriminant analysis (hyperplanes), Linear/polynomial regression, Bayesian regression)
5. Non-linear models for classification and regression (Support Vector Machines and kernel methods, Neural networks)
6. Combining models (Ensemble methods like boosting and bagging, Tree-based models)
7. Unsupervised learning methods for clustering and dimensionality reduction (E.g., hard/soft k-means clustering, Principal Component Analysis (PCA))
8. Select advanced topics based on time available (E.g., a subset of computational learning theory, algorithms for sequential data (Hidden Markov Models HMM), or graph-structured data; probabilistic graphical models).

Recall: A big question is:

WHY???

learn so many models/methods, if only one of the methods is all you hear about everywhere?

On a different context, always ask WHY you are taking this course? I hope it is not only campus placement opportunities or hype around ML, but also a general interest in understanding how you can take this field forward with your own creativity and internal ethical compass.

Parting thoughts...

Fun starts when you can take what you've learned from this course and apply it 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; different paradigms, etc.
- Understanding conceptual/mathematical foundations of different methods also helps towards above goal.

Achieving machine intelligence (PRML) is hard, but we've

- a constructive argument for many specific problems, as seen in this course, and
- an existence proof for the general problem!

Real fun is in taking this field forward with your own creativity and internal ethical compass!

Thank you...

...for your interest in the course, your attention and your feedbacks!