

Contents

- Lecture 1: Overview of Data Science
- Lecture 2: Data crawling and preprocessing
- Lecture 3: Data cleaning and integration
- Lecture 4: Exploratory data analysis
- Lecture 5: Data visualization
- Lecture 6: Multivariate data visualization
- **Lecture 7: Machine learning**
- Lecture 8: Big data analysis
- Lecture 9: Capstone Project guidance
- Lecture 10+11: Text, image, graph analysis
- Lecture 12: Evaluation of analysis results

Linear regression: introduction

- **Regression problem:** learn an unknown function $y = f(x)$ from a given training data $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_M, y_M)\}$ such that $y_i \approx f(x_i)$ for every i
 - Each observation of x is represented by a vector in an n -dimensional space, e.g., $x_i = (x_{i1}, x_{i2}, \dots, x_{in})^T$. Each dimension represents an attribute/feature/variante.
 - Bold characters denote vectors.
- **Linear model:** if $f(x)$ is assumed to be of linear form

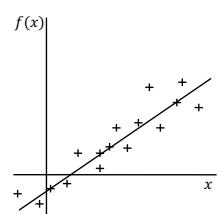
$$f(x) = w_0 + w_1x_1 + \dots + w_nx_n$$
 - w_0, w_1, \dots, w_n are the regression coefficients/weights. w_0 sometimes is called "bias".
- **Note:** learning a linear function is equivalent to learning the coefficient vector $w = (w_0, w_1, \dots, w_n)^T$.

Linear regression

Linear regression: example

- What is the best function?

x	y
0.13	-0.91
1.02	-0.17
3.17	1.61
-2.76	-3.31
1.44	0.18
5.28	3.36
-1.74	-2.46
7.93	5.56
...	...

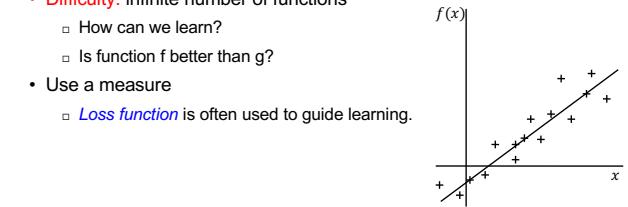


Prediction

- For each observation $\mathbf{x} = (x_1, x_2, \dots, x_n)^T$
 - The **true output**: c_x (but unknown for future data)
 - Prediction** by our system:
 $y_x = w_0 + w_1x_1 + \dots + w_nx_n$
 - We often expect $y_x \approx c_x$.
- Prediction for a future observation $\mathbf{z} = (z_1, z_2, \dots, z_n)^T$
 - Use the learned function to make prediction
 $f(\mathbf{z}) = w_0 + w_1z_1 + \dots + w_nz_n$

Learning a regression function

- Learning goal:** learn a function f^* such that its prediction in the future is the best.
 - Its generalization is the best.
- Difficulty:** infinite number of functions
 - How can we learn?
 - Is function f better than g ?
- Use a measure
 - Loss function** is often used to guide learning.



Loss function (hàm lỗi)

- Definition:**
 - The **error/loss** of the prediction for an observation $\mathbf{x} = (x_1, x_2, \dots, x_n)^T$
 $r(\mathbf{x}) = [c_x - f(\mathbf{x})]^2 = (c_x - w_0 - w_1x_1 - \dots - w_nx_n)^2$
 - The **expected loss/risk** of f over the whole space:
 $E_x[r(\mathbf{x})] = E_x[c_x - f(\mathbf{x})]^2$ (E_x is the expectation over x)
- The goal of learning is to find f^* that minimizes the expected loss:

$$f^* = \arg \min_{f \in H} E_x[r(\mathbf{x})]$$
 - H is the space of functions of linear form.
 - But, we cannot work directly with this problem during the learning phase. (why?)

Empirical loss

- We can only observe a set of training data $D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_M, y_M)\}$, and have to learn f from D.
- Residual sum of squares:

$$RSS(f) = \sum_{i=1}^M (y_i - f(\mathbf{x}_i))^2 = \sum_{i=1}^M (y_i - w_0 - w_1x_{i1} - \dots - w_nx_{in})^2$$
 - $\frac{1}{M} RSS(f)$ is an approximation to $E_x[r(\mathbf{x})]$, and is often known as **Empirical loss/risk** (lỗi thực nghiệm).
- $\left| \frac{1}{M} RSS(f) - E_x[r(\mathbf{x})] \right|$ is often known as **generalization error** of f. (lỗi tổng quát hoá)
- Many learning algorithms base on this RSS and its variants.

Methods: ordinary least squares (OLS)

- Given D, we find f^* that minimizes RSS:

$$f^* = \arg \min_{f \in H} RSS(f)$$

$$\Leftrightarrow \mathbf{w}^* = \arg \min_{\mathbf{w}} \sum_{i=1}^M (y_i - w_0 - w_1x_{i1} - \dots - w_nx_{in})^2 \quad (1)$$
- This method is often known as **ordinary least squares (OLS, binh phuong tot thi phieu)**.
- Find \mathbf{w}^* by taking the gradient of RSS and the solving the equation $RSS=0$. We have:

$$\mathbf{w}^* = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y}$$
 - Where \mathbf{A} is the data matrix of size $M \times (n+1)$, whose the i^{th} row is $\mathbf{A}_i = (1, x_{i1}, x_{i2}, \dots, x_{in})$; \mathbf{B}^{-1} is the inversion of matrix \mathbf{B} ; $\mathbf{y} = (y_1, y_2, \dots, y_M)^T$.
 - Note: we assume that $\mathbf{A}^T \mathbf{A}$ is invertible.

Methods: OLS

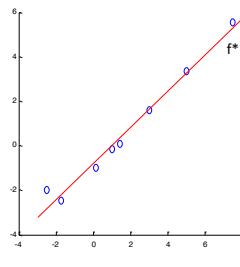
- Input: $D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_M, y_M)\}$
- Output: \mathbf{w}^*
- Learning: compute

$$\mathbf{w}^* = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y}$$
 - Where \mathbf{A} is the data matrix of size $M \times (n+1)$, whose the i^{th} row is $\mathbf{A}_i = (1, x_{i1}, x_{i2}, \dots, x_{in})$; \mathbf{B}^{-1} is the inversion of matrix \mathbf{B} ; $\mathbf{y} = (y_1, y_2, \dots, y_M)^T$.
 - Note: we assume that $\mathbf{A}^T \mathbf{A}$ is invertible.
- Prediction for a new \mathbf{x} :

$$y_x = w_0^* + w_1^*x_1 + \dots + w_n^*x_n$$

Methods: OLS example

x	y
0.13	-1
1.02	-0.17
3	1.61
-2.5	-2
1.44	0.1
5	3.36
-1.74	-2.46
7.5	5.56



Methods: limitations of OLS

- OLS cannot work if $\mathbf{A}^T\mathbf{A}$ is not invertible
 - If some columns (attributes/features) of \mathbf{A} are dependent, then \mathbf{A} will be singular and therefore $\mathbf{A}^T\mathbf{A}$ is not invertible.
- OLS requires considerable computation due to the need of computing a matrix inversion.
 - Intractable for the very high dimensional problems.
- OLS very likely tends to overfitting, because the learning phase just focuses on minimizing errors on the training data.

Methods: Ridge regression (1)

- Given $\mathbf{D} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_M, y_M)\}$, we solve for:

$$f^* = \arg \min_{f \in H} \text{RSS}(f) + \lambda \|\mathbf{w}\|_2^2$$

$$\Leftrightarrow \mathbf{w}^* = \arg \min_{\mathbf{w}} \sum_{i=1}^M (y_i - \mathbf{A}_i \mathbf{w})^2 + \lambda \sum_{j=0}^n w_j^2 \quad (2)$$

Where $\mathbf{A}_i = (1, x_{i1}, x_{i2}, \dots, x_{in})$ is composed from \mathbf{x}_i , and λ is a regularization constant ($\lambda > 0$). $\|\mathbf{w}\|_2$ is the L^2 norm.



Tikhonov,
smoothing an ill-
posed problem



Zaremba, model
complexity
minimization



Bayes: priors
over parameters



Andrew Ng: need no
maths, but it prevents
overfitting!

Methods: Ridge regression (2)

- Problem (2) is equivalent to the following:

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \sum_{i=1}^M (y_i - \mathbf{A}_i \mathbf{w})^2 \quad (3)$$

Subject to $\sum_{j=0}^n w_j^2 \leq t$ for some constant t .

- The regularization/penalty term: $\lambda \|\mathbf{w}\|_2^2$

- Limits the magnitude/size of \mathbf{w}^* (i.e., reduces the search space for f^*).
- Helps us to trade off between the fitting of f on \mathbf{D} and its generalization on future observations.

Methods: Ridge regression (3)

- We solve for \mathbf{w}^* by taking the gradient of the objective function in (2), and then zeroing it. Therefore we obtain:

$$\mathbf{w}^* = (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I}_{n+1})^{-1} \mathbf{A}^T \mathbf{y}$$

- Where \mathbf{A} is the data matrix of size $M \times (n+1)$, whose the i^{th} row is $\mathbf{A}_i = (1, x_{i1}, x_{i2}, \dots, x_{in})$; \mathbf{B}^{-1} is the inversion of matrix \mathbf{B} ; $\mathbf{y} = (y_1, y_2, \dots, y_M)^T$; \mathbf{I}_{n+1} is the identity matrix of size $n+1$.
- Compared with OLS, Ridge can
 - Avoid the cases of singularity, unlike OLS. Hence Ridge always works.
 - Reduce overfitting.
 - But error in the training data might be greater than OLS.
- Note: the quality of Ridge depends heavily on the choice of the hyperparameter λ .

Methods: Ridge regression (4)

- Input: $\mathbf{D} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_M, y_M)\}$ and $\lambda > 0$
- Output: \mathbf{w}^*
- Learning: compute

$$\mathbf{w}^* = (\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I}_{n+1})^{-1} \mathbf{A}^T \mathbf{y}$$

- Prediction for a new \mathbf{x} :

$$y_x = w_0^* + w_1^* x_1 + \dots + w_n^* x_n$$

- Note: to avoid some negative effects of the magnitude of y on covariates \mathbf{x} , one should remove w_0 from the penalty term in (2). In this case, the solution of \mathbf{w}^* should be modified slightly.

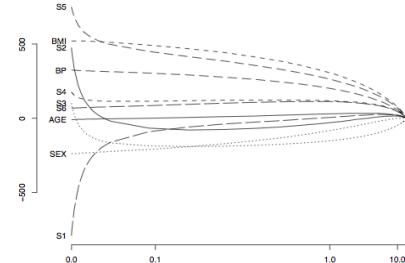
An example of using Ridge and OLS

- The training set \mathbf{D} contains 67 observations on prostate cancer, each was represented with 8 attributes. Ridge and OLS were learned from \mathbf{D} , and then predicted 30 new observations.

w	Ordinary Least Squares	Ridge
0	2.465	2.452
Icavol	0.680	0.420
Iweight	0.263	0.238
age	-0.141	-0.046
lbph	0.210	0.162
svi	0.305	0.227
Icp	-0.288	0.000
gleason	-0.021	0.040
pgg45	0.267	0.133
Test RSS	0.521	0.492

Effects of λ in Ridge regression

- $\mathbf{W}^* = (w_0, S1, S2, S3, S4, S5, S6, AGE, SEX, BMI, BP)$ changes as the regularization constant λ changes.



LASSO

- Ridge regression use L^2 norm for regularization:

$$w^* = \arg \min_{\mathbf{w}} \sum_{i=1}^M (y_i - \mathbf{A}_i \mathbf{w})^2, \text{ subject to } \sum_{j=0}^n w_j^2 \leq t \quad (3)$$

- Replacing L^2 by L^1 norm will result in LASSO:

$$w^* = \arg \min_{\mathbf{w}} \sum_{i=1}^M (y_i - \mathbf{A}_i \mathbf{w})^2$$

Subject to $\sum_{j=0}^n |w_j| \leq t$

- Equivalently:

$$w^* = \arg \min_{\mathbf{w}} \sum_{i=1}^M (y_i - \mathbf{A}_i \mathbf{w})^2 + \lambda \|\mathbf{w}\|_1 \quad (4)$$

- This problem is non-differentiable \rightarrow the training algorithm should be more complex than Ridge.

LASSO: regularization role

- The regularization types lead to different domains for \mathbf{w} .
- LASSO often produces **sparse** solutions, i.e., many components of \mathbf{w} are zero.

▫ Shinkage and feature selection at the same time

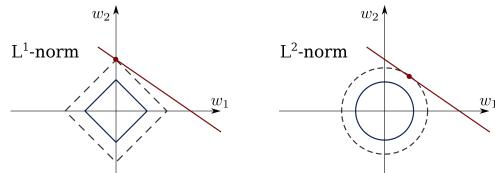


Figure by Nicoguaro - Own work, CC BY 4.0,
<https://commons.wikimedia.org/w/index.php?curid=58258966>

OLS, Ridge, and LASSO

- The training set \mathbf{D} contains 67 observations on prostate cancer, each was represented with 8 attributes. OLS, Ridge, and LASSO were trained from \mathbf{D} , and then predicted 30 new observations.

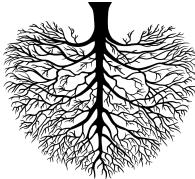
w	Ordinary Least Squares	Ridge	LASSO
0	2.465	2.452	2.468
Icavol	0.680	0.420	0.533
Iweight	0.263	0.238	0.169
age	-0.141	-0.046	
lbph	0.210	0.162	0.002
svi	0.305	0.227	0.094
Icp	-0.288	0.000	
gleason	-0.021	0.040	
pgg45	0.267	0.133	
Test RSS	0.521	0.492	0.479

Classification

Random forest

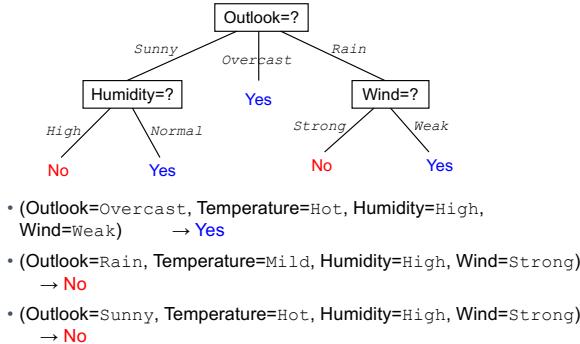
1. Decision tree

- Decision tree
 - To represent a function by using a tree.
- Each decision tree can be interpreted as a set of rules of the form: IF-THEN
- Decision trees have been used in many practical applications.



25

Examples of a decision tree (2)



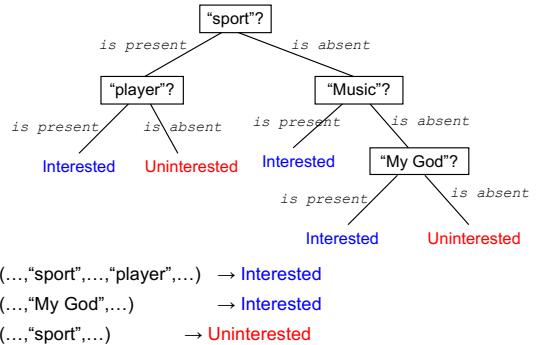
27

Tree representation (1)

- Each internal node represents an attribute for testing the incoming data.
- Each branch/subtree of a node corresponds to a value of the attribute of that node.
- Each leaf node represents a class label.
- Once a tree has been learned, we can predict the label for a new instance by using its attributes to travel from the root down to a leaf.
 - The label of the leaf will be used to assign to the new instance.

29

Examples of a decision tree (1)

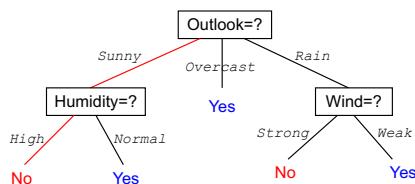


Classification problem

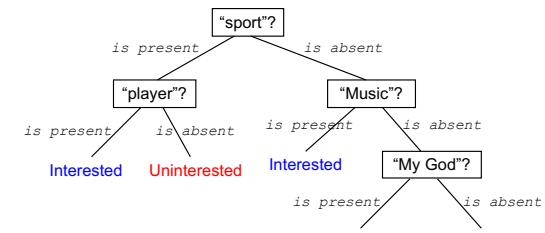
- Data representation:
 - Each observation is represented by n attributes/features, e.g., $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{in})^T$.
 - Each attribute is nominal/categorical, i.e., represents names, labels or categories, e.g., $x_{i1} \in \{high, normal\}$, $x_{i2} \in \{male, female, other\}$
- There is a set C of predefined labels.
- We have to learn a function from a training dataset:
 $D = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_M, y_M)\}$

Tree representation (2)

- Each path from the root to a leaf is a conjunction/AND of the attribute tests.
- A decision tree itself is a disjunction/OR of those conjunctions.



Representation by a disjunction



$[(\text{"sport"} \text{ is present}) \wedge (\text{"player"} \text{ is present})] \vee$
 $[(\text{"sport"} \text{ is absent}) \wedge (\text{"Music"} \text{ is present})] \vee$
 $[(\text{"sport"} \text{ is absent}) \wedge (\text{"Music"} \text{ is absent}) \wedge (\text{"My God"} \text{ is present})]$

The ID3 algorithm

ID3_alg(*Training_Set*, *Class_Labels*, *Attributes*)

Generate the Root of the tree

If all of *Training_Set* belong to class *c*, then Return Root as leaf with label *c*

If *Attributes* is empty, then Return Root as leaf with label *c* = **Majority_Class_Label**(*Training_Set*)

A \leftarrow a set of *Attributes* that are best discriminative for *Training_Set*

Let *A* be the test attributes of Root

For each value *v* of *A*

 Generate a branch of Root which corresponds with *v*.

 Determine *Training_Set_v* = {*x* in *Training_Set* | *x*_{*A*} = *v*}

 If (*Training_Set_v* is empty) Then

 Generate a leaf with class label *c* = **Majority_Class_Label**(*Training_Set*)

 Else

 Generate a subtree by **ID3_alg**(*Training_Set_v*, *Class_Labels*, *Attributes \{A\}*)

Return Root

Information gain: entropy

- Entropy measures the impurity/inhomogeneity of a set.

- Entropy of a set *S* with *c* classes can be defined as:

$$\text{Entropy}(S) = - \sum_{i=1}^c p_i \log_2 p_i$$

Where *p_i* is the proportion of instances with class label *i* in *S*; and $0 \cdot \log_0 0 = 0$ as a convention; $p_1 + p_2 + \dots + p_c = 1$

- For 2 classes: $\text{entropy}(S) = -p_1 \log_2 p_1 - p_2 \log_2 p_2$

- Meanings of entropy in Information Theory:

Entropy shows the *number of bits on average* to encode a class of *S*.

Entropy of a message measures the *average amount of information contained in that message*.

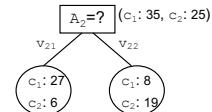
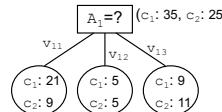
Entropy of a random variable *x* measures the *unpredictability* of *x*.

2. Learning a decision tree by ID3

- ID3 (Iterative Dichotomiser 3) is a greedy algorithm which was proposed by Ross Quinlan in 1986.
- It uses the top-down scheme.
- At each node *N*, select a **test attribute *A*** which can help us best do classification for the data in *N*.
 - Generate a branch for each value of *A*, and then separate the data into its branches accordingly.
- Grow the tree until:
 - It classifies correctly all the training data; or
 - All the attributes are used.
- Note:** each attribute can only appear at most once in any path of the tree.

How to choose the test attributes?

- At each node, **how can we choose a set of test attributes?**
 - These attributes should be **discriminative**, i.e., can help us classify well the data inside that node.
- How to know an attribute to be discriminative?**
- Ex: assuming 2 classes in the data, which of *A*₁ and *A*₂ should be selected as the test attribute?



Information gain can help.

Information gain: entropy example

- S* consists of 14 examples for which 9 belong to class *c*₁ and 5 belong to class *c*₂.

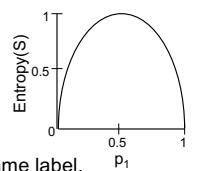
- So the entropy of *S* is:

$$\begin{aligned} \text{Entropy}(S) &= -(9/14) \cdot \log_2(9/14) - (5/14) \cdot \log_2(5/14) \\ &\approx 0.94 \end{aligned}$$

- Entropy = 0 if all examples in *S* have the same label.

- Entropy = 1 if the two classes in *S* are equal in size.

- Otherwise, entropy will always belong to (0, 1).



Information gain

- **Information gain** of an attribute in S:

▫ Measures the reduction of entropy if we divide S into subsets according to that attribute.

- Information gain of attribute A in S is defined as:

$$Gain(S, A) = Entropy(S) - \sum_{v \in Values(A)} \frac{|S_v|}{|S|} Entropy(S_v)$$

▫ Where $Values(A)$ is the set of all values of A, and $S_v = \{x | x \in S, \text{ and } x_A = v\}$

- The **second term** in $Gain(S, A)$ measures the *information remained* when S is divided into subsets according to the values of A.

- **Meaning of $Gain(S, A)$:** *the average amount of information is lost when dividing S according to A.*

Information gain: example (2)

- What is $Gain(S, Wind)$?

- Wind has two values: Strong & Weak

- $S = \{9 \text{ examples with label Yes}, 5 \text{ examples with label No}\}$

- $S_{Weak} = \{6 \text{ examples with label Yes and 2 examples with label No, having Wind=Weak}\}$

- $S_{Strong} = \{3 \text{ examples with label Yes, 3 examples with label No, having Wind=Strong}\}$

- So:

$$\begin{aligned} Gain(S, Wind) &= Entropy(S) - \sum_{v \in \{Strong, Weak\}} \frac{|S_v|}{|S|} Entropy(S_v) \\ &= Entropy(S) - \frac{8}{14} Entropy(S_{Weak}) - \frac{6}{14} Entropy(S_{Strong}) \\ &= 0.94 - \frac{8}{14} * 0.81 - \frac{6}{14} * 1 = 0.048 \end{aligned}$$

Information gain: example (1)

- A set S of observations about a person playing tennis.

Day	Outlook	Temperature	Humidity	Wind	Play Tennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

ID3: example (2)

- At Node1, which one of {Temperature, Humidity, Wind} should be the test attribute?

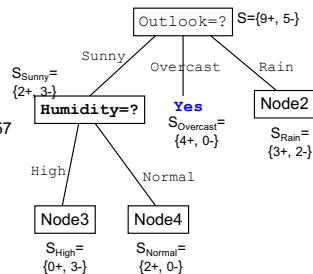
▫ Note: Outlook is left out

▫ $Gain(S_{Sunny}, Wind) = \dots = 0.019$

▫ $Gain(S_{Sunny}, Temperature) = \dots = 0.57$

▫ $Gain(S_{Sunny}, \text{Humidity}) = \dots = 0.97$

- So, Humidity is selected to divide Node1.



ID3: searching scheme (1)

- ID3 searches for a tree that fits well with the training data.

▫ By growing the tree gradually.

- **Information Gain decides the search direction of ID3.**

- ID3 just searches for **only one tree**.

- ID3 never backtracks, as a consequence:

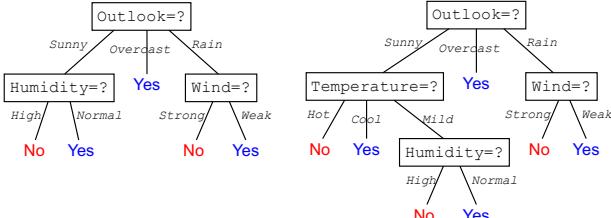
▫ It can find a **local optimal solution/tree**.

▫ Once an attribute has been selected, ID3 never rethinks of this choice.

ID3: searching scheme (2)

- For a training dataset, there might be many trees that fit well with it.

- Which tree will be selected by ID3?



ID3: searching scheme (3)

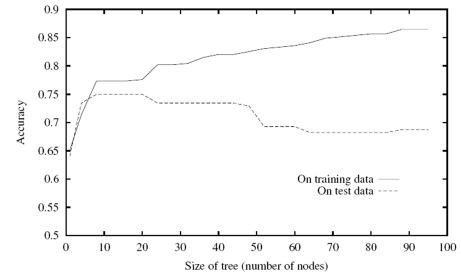
- ID3 selects the first tree that fits the training data,
 - Because it never reconsiders its choices when growing a tree.
- So, the searching scheme of ID3:
 - Prefers simple trees.
 - Prefers trees in which the attributes with higher information gain will be placed closer to the roots.

3. Some issues of ID3

- The learnt trees may overfit the training data.
- How to work with real attributes?
 - Many applications have real inputs.
- Is there any better measure than information gain?
- How to deal with missing values?
 - Missing-value is an inherent problem in many practical applications.
- How to enclose the cost of attributes in ID3?

Overfitting in ID3 (2)

- An example: continuing to grow the tree can improve the accuracy on the training data, but perform badly on the test data.



Overfitting: solutions

- 2 solutions:
 - Stop learning early:** prevent the tree before it fits the training data perfectly.
 - Prune the full tree:** grow the tree to its full size, and then post prune the tree.
- It is hard to decide when to stop learning.
- Post-pruning the tree empirically results in better performance. But
 - How to decide the good size of a tree?
 - When to stop pruning?
- We can use a validation set to do pruning, such as, *reduced-error pruning*, and *rule-post pruning*.

ID3: attribute selection

- Information gain:
 - Prefers the attribute that has more unique values.
 - Attributes with more unique values will be placed closer to the root than the other attribute.
- We can use some other measures, such as **Gain Ratio**

$$GainRatio(S, A) = \frac{Gain(S, A)}{SplitInformation(S, A)},$$

$$SplitInformation(S, A) = - \sum_{v \in Values(A)} \frac{|S_v|}{|S|} \log_2 \frac{|S_v|}{|S|}$$

ID3: missing or real values

• How to work with real attributes?

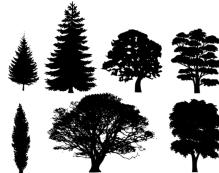
- Real attributes/features are popular in practice.
- One way is to *discretization*, i.e., transforming a real attribute into a discrete one by dividing the domain of that attribute into a set of intervals.
Ex: $[0, 1] \rightarrow \{[0, 0.25); [0.25, 0.5); [0.5, 0.75); [0.75, 1]\}$

• How to deal with missing values?

- Missing values are inherent in practical applications.
- An observation \mathbf{x} may not have a value x_A .
- *Solution 1:* fill in x_A as the most popular value of A in the training data.
- *Solution 2:* fill in x_A as the most popular value of A in the training data which belong to the same class with \mathbf{x} .

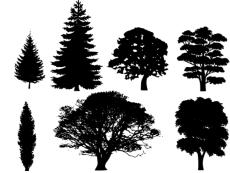
5. Random forests

- RF currently is one of the most popular and accurate methods [Fernández-Delgado et al., 2014]
 - It is also very general.
- RF can be implemented easily and efficiently.
- It can work with problems of **very high dimensions**, without overfitting ☺
- However, little is known about its theoretical properties ☺



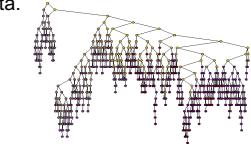
5. Random forests

- Random forests (RF) is a method by Leo Breiman (2001) for both classification and regression.
- **Main idea:** prediction is based on combination of many decision trees, by *taking the average of all individual predictions*.
- Each tree in RF is simple but random.
- Each tree is grown differently, depending on the choices of the attributes and training data.



5. RF: three basic ingredients

- **Randomization and no pruning:**
 - For each tree and at each node, we select randomly a subset of attributes.
 - Find the best split, and then grow appropriate subtrees.
 - Every tree will be grown to its largest size without pruning.
- **Combination:** each prediction later is made by taking the average of all predictions of individual trees.
- **Bagging:** the training set for each tree is generated by sampling (with replacement) from the original data.



5. RF: algorithm

- **Input:** training data D
- **Learning:** grow K trees as follows
 - Generate a training set D_i by sampling with replacement from D.
 - Learn the i^{th} tree from D_i :
 - At each node:
 - Select randomly a subset S of attributes.
 - Split the node into subtrees according to S.
 - Grow this tree upto its largest size without pruning.
 - **Prediction:** taking the average of all predictions from the individual trees.

5. RF: practical performance

- RF is extensively compared with other methods
 - By Fernández-Delgado et al. (2014).
 - Using 55 different problems.
 - Using average accuracy (μ^P) as a measure.

No.	Classifier	μ^P	No.	Classifier	μ^P
1	rf.t	91.1	11	Bagging.LibSVM.w	89.9
2	parRF.t	91.1	12	RandomCommittee.w	89.9
3	svm.C	90.7	13	Bagging.RandomTree.w	89.8
4	RRF.t	90.6	14	MultiBoostAB.RandomTree.w	89.8
5	RRFglobal.t	90.6	15	MultiBoostAB.LibSVM.w	89.8
6	LibSVM.w	90.6	16	MultiBoostAB.PART.w	89.7
7	RotationForest.w	90.5	17	Bagging.PART.w	89.7
8	C5.0.t	90.5	18	AdaBoostM1.J48.w	89.5
9	rforest.R	90.3	19	Bagging.REPTree.w	89.5
10	treebag.t	90.2	20	MultiBoostAB.J48.w	89.4



VIỆN CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG
SCHOOL OF INFORMATION AND COMMUNICATION TECHNOLOGY

Thank you
for your
attentions!

soict.hust.edu.vn/ fb.com/groups/soict

