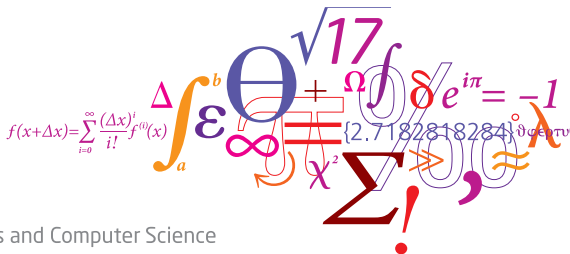


02450: Introduction to Machine Learning and Data Mining

Decision trees and linear regression



Reading Material

Reading material:

C7, C8

Feedback Groups of the day:

- Albert Juhl, Mathias Henriksen
- Marianne Helenius, Niklas Refsgaard, Martin Haubro
- Martin Petersson, Christoffer Jensen
- Mads Okholm Bjørn, Rasmus Liebst, Johan Bloch Madsen
- Line Maj Thomsen, Johan Lassen
- Ramiro Mata, Bianca Burger, Marsela Fallah
- Martin Simon, Péter Semság
- Thomas Masquart, Julien Hoareau

Tue Herlau, Mikkel N. Schmidt and Morten Mørup

Introduction to Machine Learning and
Data Mining

Course notes fall 2016, version 1

August 29, 2016

Technical University of Denmark

Lecture Schedule

1 Introduction

30 August: C1

Data: Feature extraction, and visualization

2 Data and feature extraction

6 September: C2, C3

3 Measures of similarity and summary statistics

13 September: C4

4 Data Visualization and probability

20 September: C5, C6

Supervised learning: Classification and regression

5 Decision trees and linear regression

27 September: C7, C8 (Project 1 due before 13:00)

6 Overfitting and performance evaluation

4 October: C9

7 Nearest Neighbor, Bayes and Naive Bayes

11 October: TBA

8 Artificial Neural Networks and Bias/Variance

25 October: TBA

9 AUC, ensemble methods and multi-class classifiers

1 November: TBA

Unsupervised learning: Clustering and density estimation

10 K-means and hierarchical clustering

8 November: TBA (Project 2 due before 13:00)

11 Mixture models and association mining

15 November: TBA

12 Density estimation and anomaly detection

22 November: TBA

Recap

13 Recap and discussion of the exam

29 November: TBA (Project 3 due before 13:00)

Probabilities (revisited from last week)

- Basic rules of probability

- Sum rule

$$p(x) = \sum_y p(x, y)$$

- Product rule

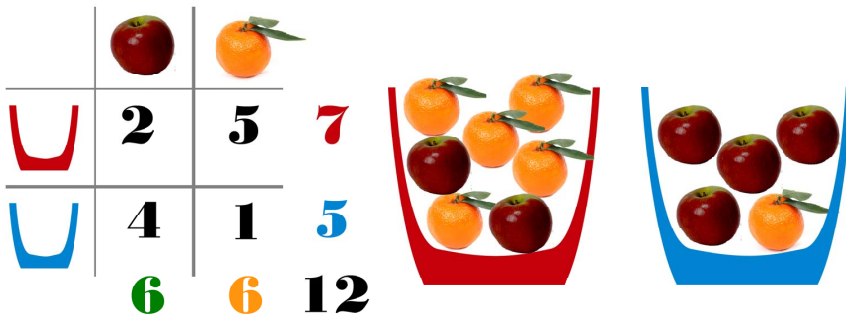
$$p(x, y) = p(x|y)p(y)$$

- Bayes' rule

$$p(x|y) = \frac{p(y|x)p(x)}{p(y)}$$

Probabilities





- What is the probability of an **orange** if the bowl is **red**?
- What is the probability of the **red** bowl if the fruit is **orange**?



Apple taken from: https://upload.wikimedia.org/wikipedia/commons/3/32/Dark_apple.png
Orange (clementine) taken from: https://commons.wikimedia.org/wiki/File:Clementine_orange.jpg

Probabilities

- What is the probability of an **orange** if the bowl is **red**?
- What is the probability of the **red** bowl if the fruit is **orange**?

			
	2	5	7
	4	1	5
	6	6	12

$$p(o|r) = \frac{p(r, o)}{p(r)} = \frac{5/12}{7/12} = 5/7$$

$$\begin{aligned}
 p(r|o) &= \frac{p(r, o)}{p(o)} = \frac{5/12}{6/12} = 5/6 \\
 &= \frac{p(o|r)p(r)}{p(o)} \\
 &= \frac{5/7 \cdot 7/12}{6/12} = 5/6
 \end{aligned}$$

Apple taken from: https://upload.wikimedia.org/wikipedia/commons/3/32/Dark_apple.png

Orange (clementine) taken from: https://commons.wikimedia.org/wiki/File:Clementine_orange.jpg



The news paper "Economist"

News media agency Reuters Bureau sends stories to the Economist:



- 80% of the news stories from Reuters are positive and 20% of the news stories are negative.
- 90% of the negative news stories are published in the Economist while only 5% of the positive stories are published.

Consider a story from Reuters. What is the probability it is positive given it is published in the Economist?

Hints:

- *Sum Rule*

$$p(x) = \sum_y p(x, y)$$

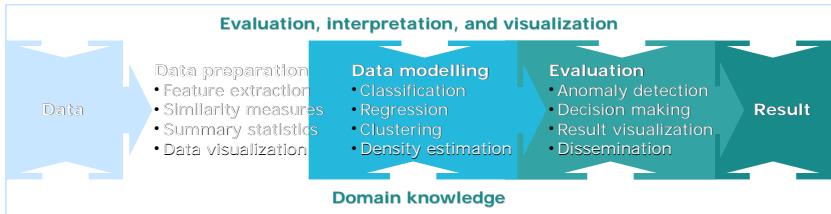
- *Product rule*

$$p(x, y) = p(x|y)p(y)$$

- *Bayes theorem*

$$p(x|y) = \frac{p(y|x)p(x)}{p(y)}$$

Data modeling framework



After today you should be able to:

Explain what supervised learning is

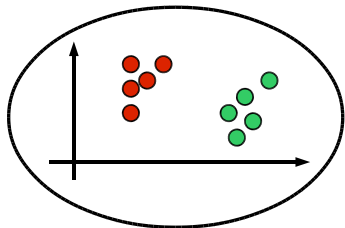
Explain the difference between classification and regression

Be able to evaluate classifiers in terms of the confusion matrix, error rate and accuracy

Understand the principles behind decision trees and Hunt's algorithm

Apply and interpret decision trees, linear regression and logistic regression

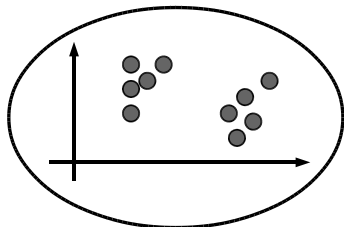
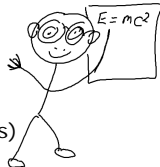
Supervised and Unsupervised learning



Supervised Learning

Input data \mathbf{x}_n and output y_n

(Generalize from known examples)



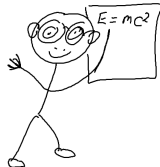
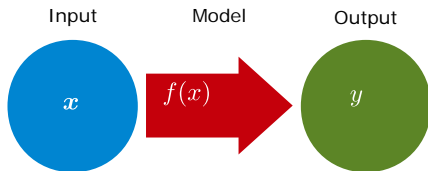
Unsupervised Learning

Input data \mathbf{x}_n alone

(Exploratory analysis)



Supervised learning



- **Data**

- Inputs and outputs (*this is what we are given*)

- **Model**

$$\{\mathbf{x}_n, y_n\}_{n=1}^N$$

- Function that maps inputs to outputs (*what we are trying to determine*)

$$f(\mathbf{x})$$

- **Cost function**

- Dissimilarity measure between observation and prediction (*how we tell if a model is good or bad*)

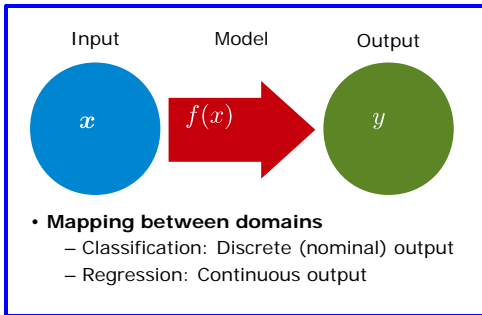
$$d(y, f(\mathbf{x}))$$

- **Types of supervised learning**

- Regression: Continuous output \mathbf{y}
- Classification: Discrete output \mathbf{y}



Give an example of a classification and a regression problem and explain what the model $f(x)$ can be used for.



Classification

- **Definition:** Learning a function that maps a data object to a discrete class
- **Why classify?**
 - Descriptive modeling
 - Explain / understand the relation between attributes and class
 - Predictive modeling
 - Predict the class of a new data object

Confusion matrix

- Visualization of actual versus predicted class labels

- **Accuracy**

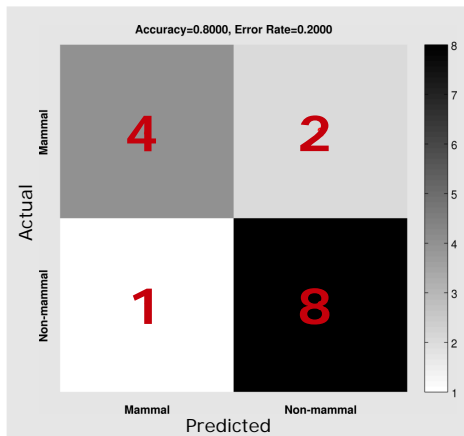
(Number of correctly predicted observations divided by the total number of observations)

$$\frac{4 + 8}{4 + 2 + 1 + 8} = 80\%$$

- **Error rate**

(Number of in-correctly predicted observations divided by the total number of observations)

$$\frac{2 + 1}{4 + 2 + 1 + 8} = 20\%$$



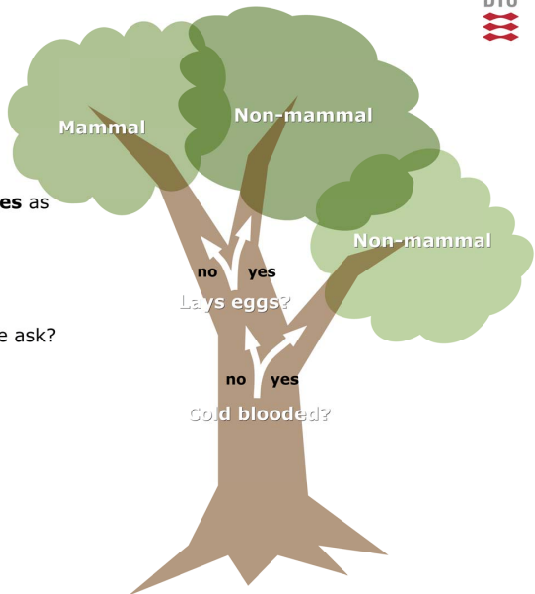
Decision trees

- Remember the game “20 questions to the professor”? (see also www.20q.new)

- Q1. Is it an Animal? Yes.
- Q2. Can you hold it? No.
- Q3. Does it live in groups (gregarious)? Yes.
- Q4. Are there many different sorts of it? No.
- Q5. Can it jump? Yes.
- Q6. Does it eat seeds? No.
- Q7. Is it white? Sometimes.
- Q8. Is it black and white? No.
- Q9. Does it have paws? Yes.
- Q10. Can you see it in a zoo? Yes.
- Q11. Does it roar? Yes.
- Q12. Is it worth a lot of money? Yes.
- Q13. Does it have spots? Yes.
- Q14. Is it multicoloured? Yes.
- Q15. Can you make money by selling it? Yes.
- Q16. Does it live in the jungle? Yes.
- Q17. I guessed that it was a leopard? Wrong.
- Q18. Does it like to play? Yes.
- Q19. I guessed that it was a cheetah? Wrong.
- Q20. I am guessing that it is a siberian tiger? Correct.

Decision trees

- Ask a series of questions until a conclusion is reached
- **Example:** Classify **vertebrates** as
 - **Mammal** or
 - **Non-mammal**
- **Learning task**
 - Which questions should we ask?



Hunts algorithm

- Assign all data objects to the root



Mammals 5:10 Non-mammals

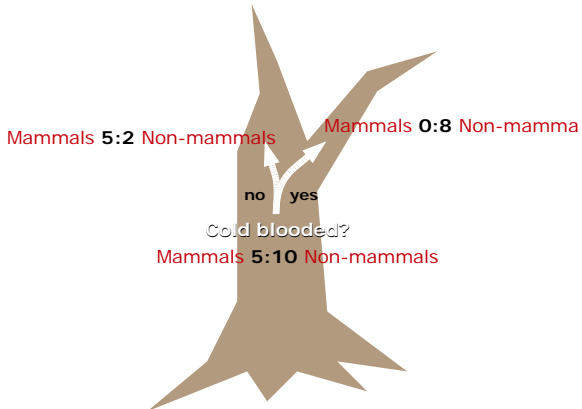
Hunts algorithm

- Select an attribute test condition
 - Find a good question to ask



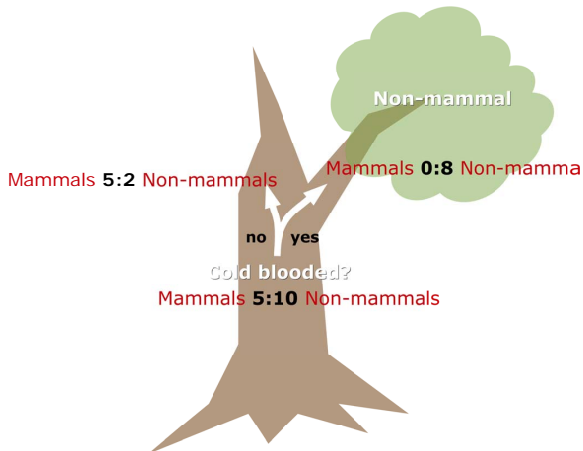
Hunt's Algorithm

- Partition the data objects into subsets according to the test condition



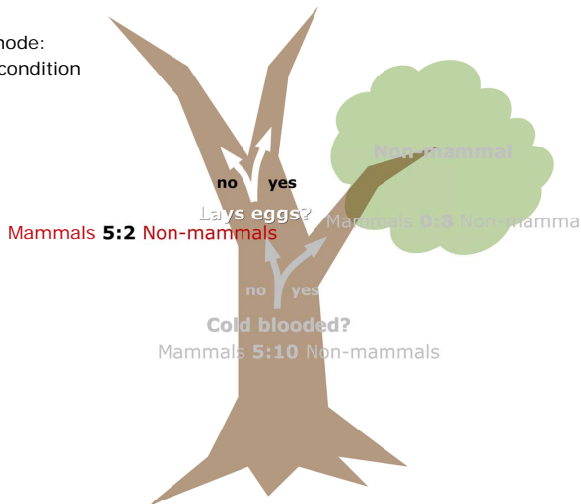
Hunts algorithm

- If all data objects belong to the same class
 - Create a leaf node



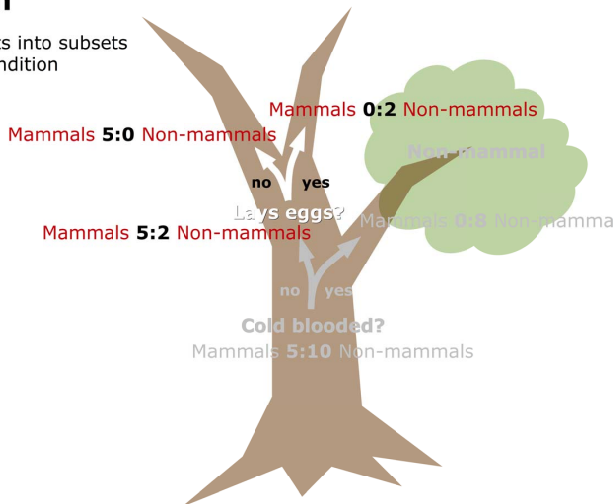
Hunts algorithm

- Repeat for each non-leave node:
 - Select an attribute test condition



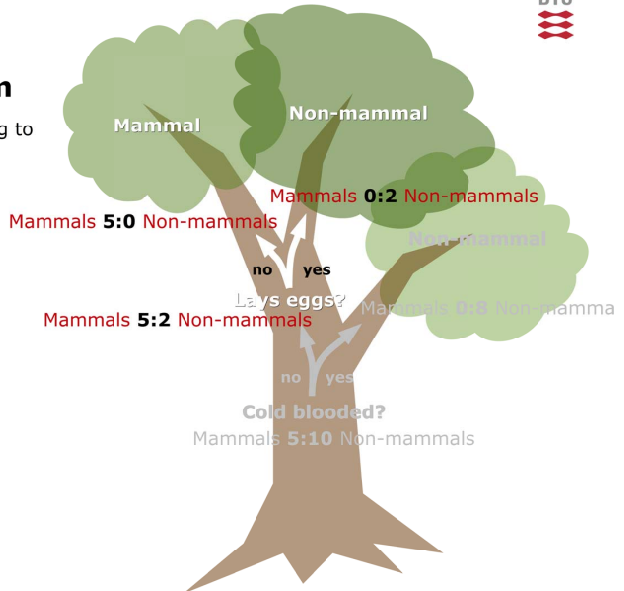
Hunts algorithm

- Partition the data objects into subsets according to the test condition



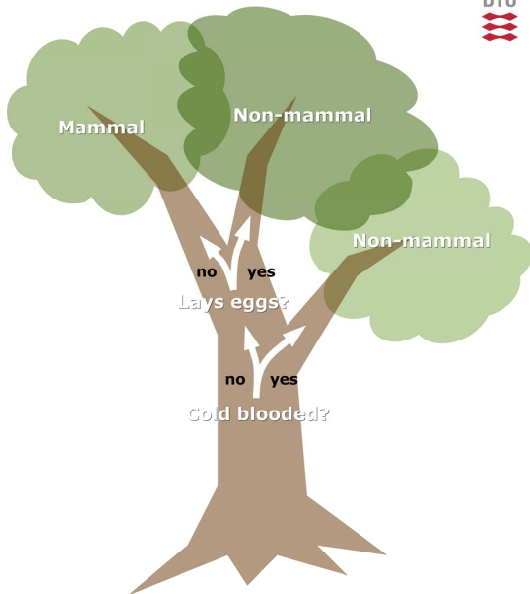
Hunts algorithm

- If all data objects belong to the same class
 - Create a leaf node



Hunts algorithm

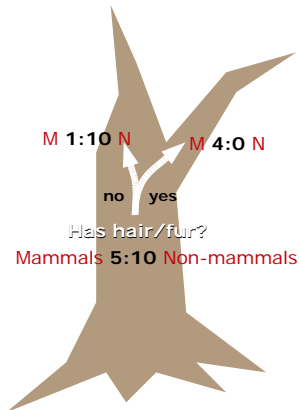
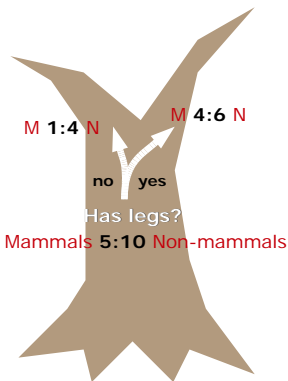
- But how do we find the **best question** at each step?





Selecting the best split

- Which of these two questions is best and why?



Selecting the best split

- Consider a large number of possible splits
- Compute a measure of impurity after the proposed split
 - For each new branch of the tree
 - Compute weighted average impurity
- Choose split that reduces impurity most



Selecting the best split: Impurity measures



- Compute the purity gain, Δ

$$Entropy(t) = - \sum_{i=0}^{c-1} p(i|t) \log_2 p(i|t)$$

$\Delta = ?$

$\Delta = ?$

$$Gini(t) = 1 - \sum_{i=0}^{c-1} (p(i|t))^2$$

$\Delta = ?$

$\Delta = ?$

$$Class. error(t) = 1 - \max_i p(i|t)$$

$\Delta = ?$

$\Delta = ?$

$p(i|t)$ Fraction of objects that belong to class i
 $N(v_j)/N$ Fraction of animals in branch v_j

$$\Delta = I(\text{parent}) - \sum_{j=1}^k \frac{N(v_j)}{N} I(v_j)$$

M 1:4 N M 4:6 N

M 1:10 N M 4:0 N





Selecting the best split: Impurity measures

- Compute the purity gain, Δ

$$\text{Entropy}(t) = - \sum_{i=0}^{c-1} p(i|t) \log_2 p(i|t)$$

$$\text{Gini}(t) = 1 - \sum_{i=0}^{c-1} (p(i|t))^2$$

$$\text{Class. error}(t) = 1 - \max_i p(i|t)$$

$p(i|t)$ Fraction of objects that belong to class i
 $N(v_j)/N$ Fraction of animals in branch v_j

$$\Delta = I(\text{parent}) - \sum_{j=1}^k \frac{N(v_j)}{N} I(v_j)$$

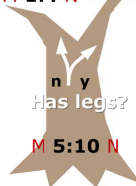
$I(\text{Parent}) = -5/15 \log(5/15) - 10/15 \log(10/15)$ $= 0.9183$ $I(\text{left}) = -1/5 \log(1/5) - 4/5 \log(4/5)$ $= 0.7219$ $I(\text{right}) = -4/10 \log(4/10) - 6/10 \log(6/10)$ $= 0.9710$ $\Delta = 0.9183 - 5/15 \cdot 0.7219 - 10/15 \cdot 0.9710$ $= 0.0303$	$I(\text{Parent}) = -5/15 \log(5/15) - 10/15 \log(10/15)$ $= 0.9183$ $I(\text{left}) = -1/11 \log(1/11) - 10/11 \log(10/11)$ $= 0.4395$ $I(\text{right}) = -4/4 \log(4/4) - 0/4 \log(0/4)$ $= 0$ $\Delta = 0.9183 - 1/15 \cdot 0.4395 - 4/15 \cdot 0$ $= 0.5960$
--	---

$I(\text{Parent}) = 1 - (5/15)^2 - (10/15)^2$ $= 0.4444$ $I(\text{left}) = 1 - (1/5)^2 - (4/5)^2$ $= 0.3200$ $I(\text{right}) = 1 - (4/10)^2 - (6/10)^2$ $= 0.4800$ $\Delta = 0.4444 - 5/15 \cdot 0.3200 - 10/15 \cdot 0.4800$ $= 0.0177$	$I(\text{Parent}) = 1 - (5/15)^2 - (10/15)^2$ $= 0.4444$ $I(\text{left}) = 1 - (1/11)^2 - (10/11)^2$ $= 0.1653$ $I(\text{right}) = 1 - (4/4)^2 - (0/4)^2$ $= 0$ $\Delta = 0.4444 - 1/15 \cdot 0.1653 - 4/15 \cdot 0$ $= 0.3232$
--	--

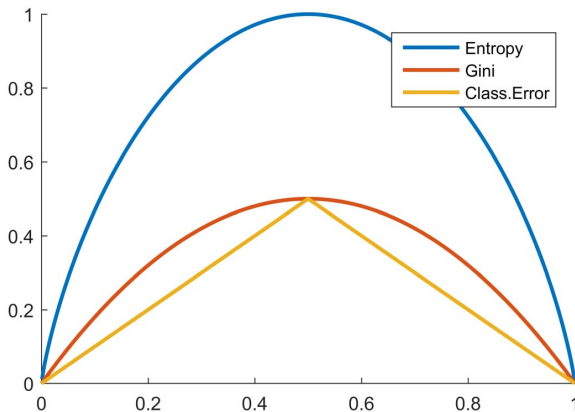
$I(\text{Parent}) = 1 - 10/15$ $= 5/15$ $I(\text{left}) = 1 - 4/5$ $= 1/5$ $I(\text{right}) = 1 - 6/10$ $= 4/10$ $\Delta = 5/15 - 5/15 \cdot 1/5 - 10/15 \cdot 4/10$ $= 0$	$I(\text{Parent}) = 1 - 10/15$ $= 5/15$ $I(\text{left}) = 1 - 10/11$ $= 1/11$ $I(\text{right}) = 1 - 4/4$ $= 0$ $\Delta = 5/15 - 11/15 \cdot 1/11 - 4/15 \cdot 0$ $= 0.2667$
---	---

M 1:4 N M 4:6 N

M 1:10 N M 4:0 N



For a two class problem



Which splits to consider



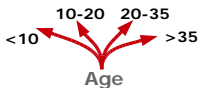
- Nominal



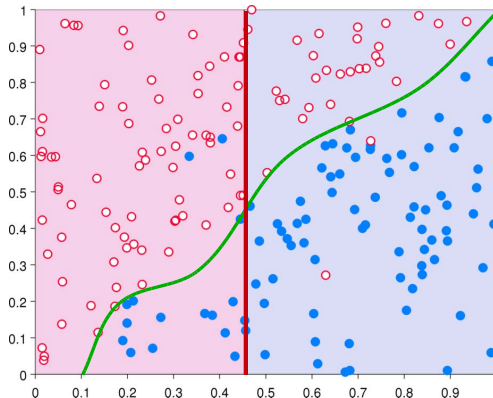
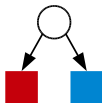
- Ordinal



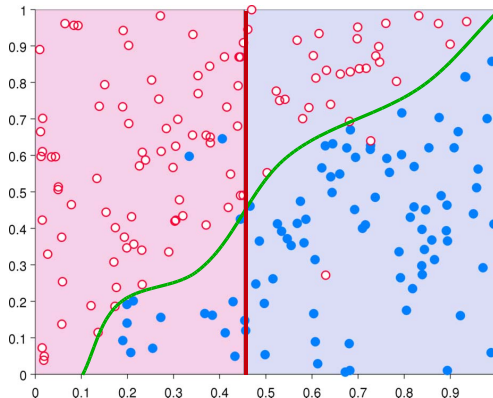
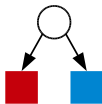
- Continuous



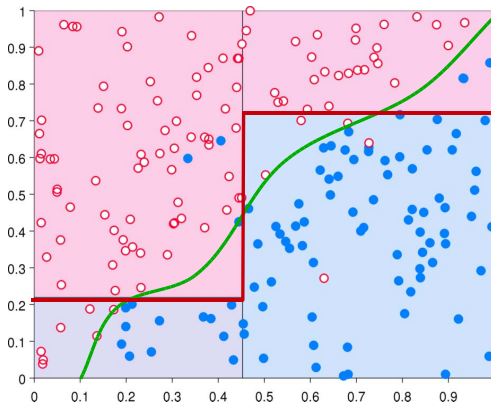
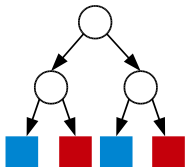
Classification Trees



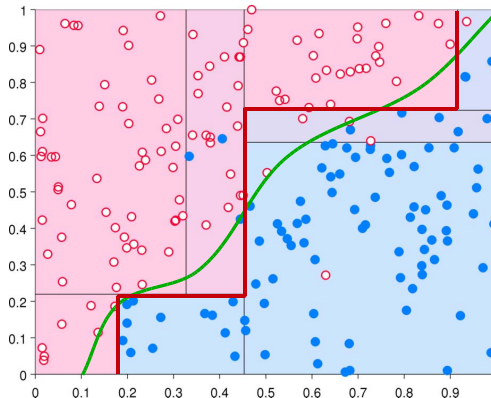
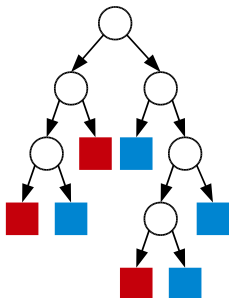
Classification Trees



Classification trees



Classification trees



All records have the same class label

The number of observations have fallen below some minimum treshhold



The iris data set

- **Three flowers**

- 50 instances of each class, 150 in total

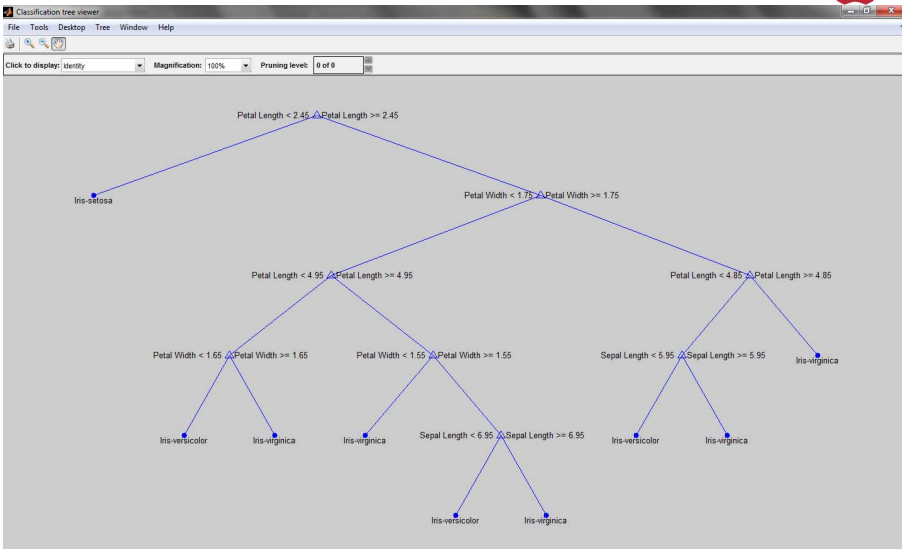
- **Attributes**

- Sepal (outermost leaves)
 - length in cm
 - width in cm
- Petal (innermost leaves)
 - length in cm
 - width in cm
- Class of flower
 - Iris Setosa
 - Iris Versicolour
 - Iris Virginica



Fower ID	Attribute			
	Sepal Length	Sepal Width	Petal Length	Petal Width
1	5.1	3.5	1.4	0.2
2	4.9	3.0	1.4	0.2
3	4.7	3.2	1.3	0.2
4	4.6	3.1	1.5	0.2
•	•	•	•	•
•	•	•	•	•
150	5.9	3.0	5.1	1.8

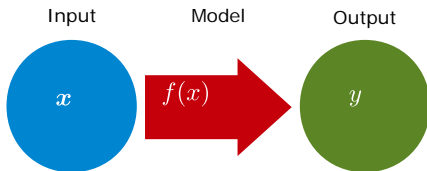
X Observation x Attribute



What would the following iris flower be classified as?

Sepal Length	Sepal Width	Petal Length	Petal Width
4.0	3.5	3.0	2.0

Supervised learning



- **Mapping between domains**

- Classification: Discrete (nominal) output
- Regression: Continuous output

Supervised learning

- **Data**

- Inputs and outputs $\{\mathbf{x}_n, y_n\}_{n=1}^N$

- **Model**

- Function that maps inputs to outputs

$$f(\mathbf{x})$$

- **Cost function**

- Dissimilarity measure between data and model

$$d(y, f(\mathbf{x}))$$

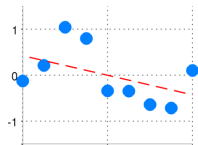
Regression

- **Definition:** Learning a function that maps a data object to a continuous-valued output
- **Why Regression?**
 - Descriptive modeling
 - Explain / understand the relation between attributes and continuous-valued output
 - Predictive modeling
 - Predict the output value of a new data object

Linear regression

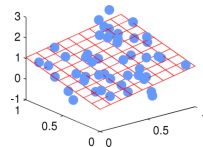
- 1-dimensional inputs

$$f(x) = w_0 + w_1x$$



- 2-dimensional inputs

$$f(\mathbf{x}) = w_0 + w_1x_1 + w_2x_2$$



- K-dimensional inputs

$$f(\mathbf{x}) = w_0 + w_1x_1 + w_2x_2 + \cdots + w_Kx_K$$

Linear regression

- K-dimensional inputs

$$f(\mathbf{x}) = w_0 + w_1x_1 + w_2x_2 + \cdots + w_Kx_K$$

- Non-linearly transformed inputs

$$f(x) = w_0 + w_1x + w_2x^2 + \cdots + w_Kx^K$$

$$f(x) = w_0 + w_1\sin(x) + w_2\cos(x)$$

Linear regression

- K-dimensional inputs

$$f(\mathbf{x}) = w_0 + w_1x_1 + w_2x_2 + \cdots + w_Kx_K$$

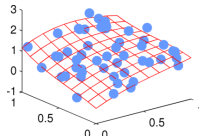
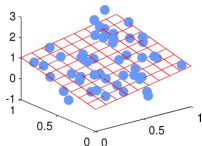
- Non-linearly transformed inputs

$$f(x) = w_0 + w_1x + w_2x^2 + \cdots + w_Kx^K$$

$$f(x) = w_0 + w_1\sin(x) + w_2\cos(x)$$

- **Example**

$$f(\mathbf{x}) = w_0 + w_1x_1 + w_2x_2$$



$$\begin{aligned} f(\mathbf{x}) = & w_0 + w_1x_1 + w_2x_2 \\ & + w_3x_1^2 + w_4x_2^2 + w_5x_1x_2 \\ & + w_6x_1^3 + w_7x_1^2x_2 + w_8x_1x_2^2 + w_9x_2^3 \end{aligned}$$

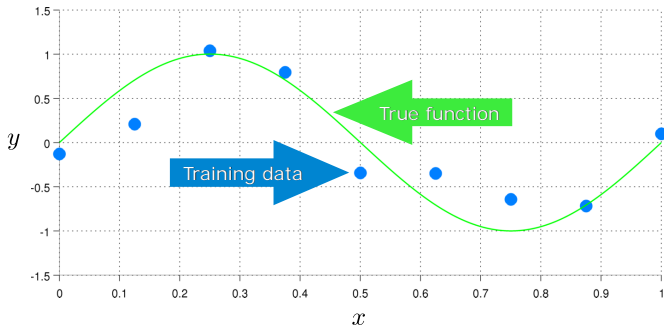
Vector notation

- The linear model can be written compactly using vector notation

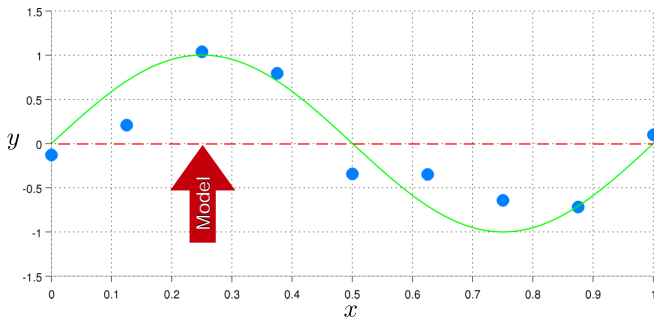
$$\begin{aligned} f(\mathbf{x}) &= w_0 + w_1 x_1 + w_2 x_2 + \cdots + w_K x_K \\ &= \sum_{k=0}^K w_k x_k = \boxed{\mathbf{x}^\top \mathbf{w}} \end{aligned}$$

– where $x_0 = 1$

Linear regression



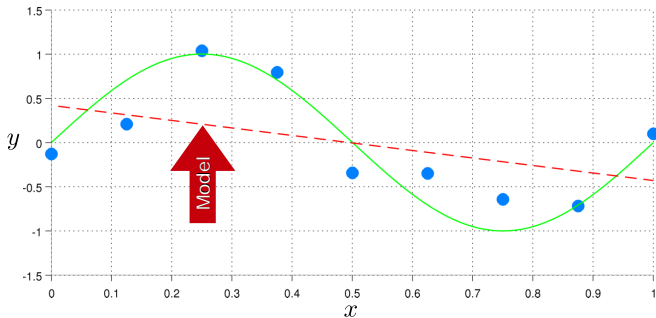
Linear regression



Model

$$f(x) = w_0$$

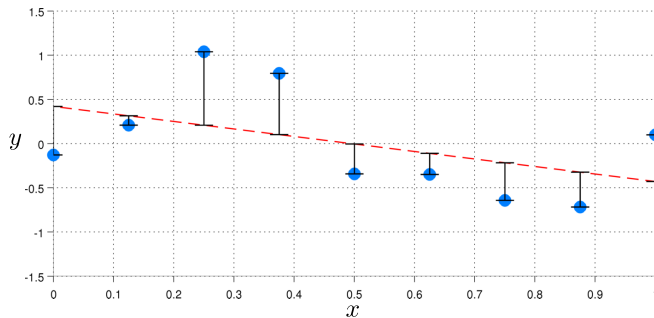
Linear regression



Model

$$f(x) = w_0 + w_1x$$

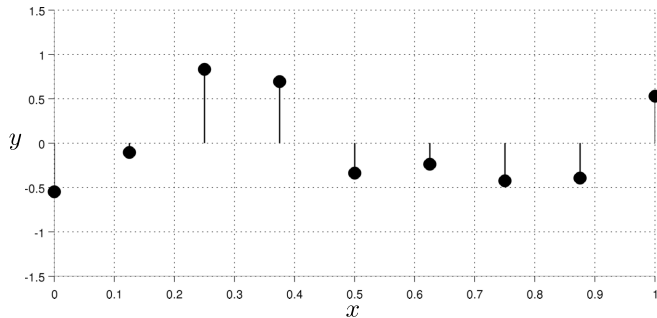
Residual error



Model

$$f(x) = w_0 + w_1 x$$

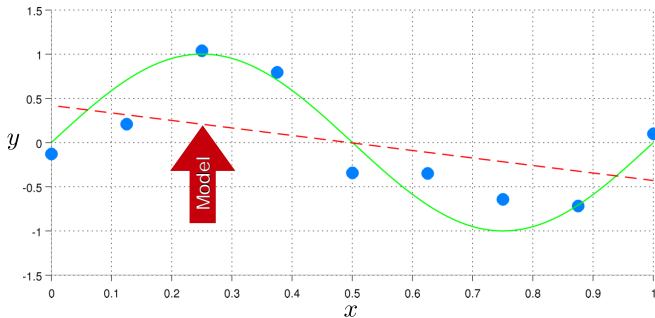
Residual error



Model

$$f(x) = w_0 + w_1x$$

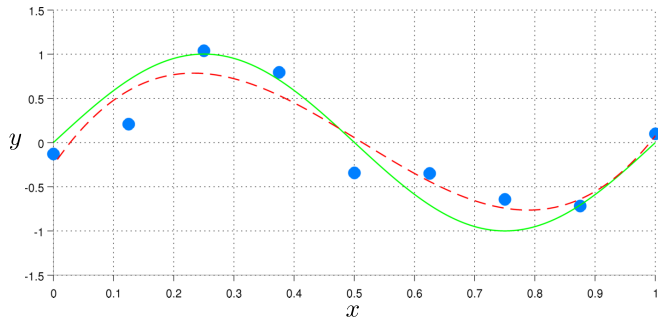
Linear regression



Model

$$f(x) = w_0 + w_1x$$

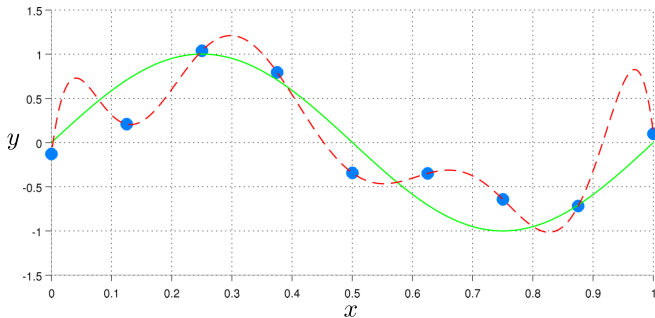
Linear regression



Model

$$f(x) = w_0 + w_1x + w_2x^2 + w_3x^3$$

Linear regression



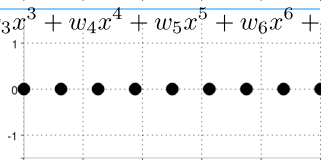
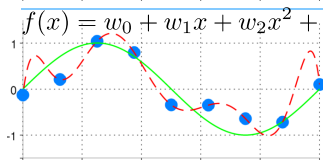
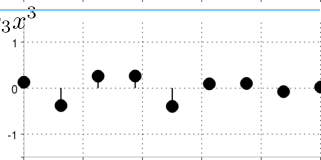
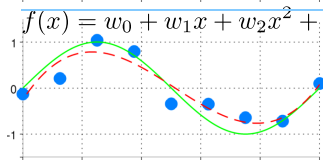
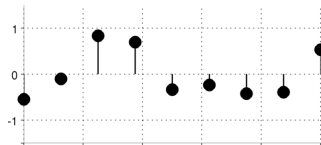
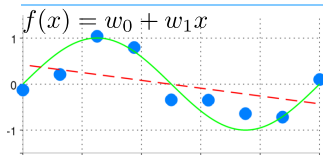
Model

$$f(x) = w_0 + w_1x + w_2x^2 + w_3x^3 + w_4x^4 + w_5x^5 + w_6x^6 + w_7x^7 + w_8x^8$$



Model order

- Which model order
 - Gives the best fit?
 - Do you think is most "correct"?



Estimating parameters

- How do we compute the parameters?

- Most simple approach: **Minimize cost function over data set**

- **Data** $\{\mathbf{x}_n, y_n\}_{n=1}^N$

- **Model** $f(\mathbf{x}) = \mathbf{x}^\top \mathbf{w}$

- **Cost function** $d(y, f(\mathbf{x}))$

- **Parameters** $\mathbf{w} = \arg \min_{\mathbf{w}} \sum_{n=1}^N d(y_n, f(\mathbf{x}_n))$

Least Squares Regression

- **Cost function:** Squared error


$$d(y, f(x)) = (y - f(x))^2$$

- **Model:** Linear regression

$$f(x) = x^\top w$$

- **Parameters**

$$w = \arg \min_w \sum_{n=1}^N d(y_n, f(x_n)) = \arg \min_w \sum_{n=1}^N (y_n - x_n^\top w)^2$$


$$E = \|y - Xw\|^2$$

$$\frac{\partial E}{\partial w} = 2(y - Xw)^\top X = 0$$

$$\Rightarrow 2y^\top X = 2w^\top X^\top X$$

$$\Rightarrow w = (X^\top X)^{-1} X^\top y$$

Logistic Regression

(for binary classification, $y \in \{0,1\}$)

- Cost function:**

negative log of the Bernoulli distribution

$$d(y, f(\mathbf{x})) = -y \log f(\mathbf{x}) - (1 - y) \log(1 - f(\mathbf{x}))$$

Model: Logistic link function

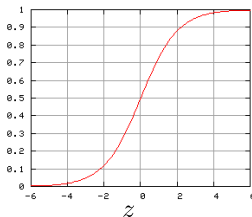
$$f(\mathbf{x}) = g_{\text{logistic}}(\mathbf{x}^\top \mathbf{w})$$

- Parameters**

$$\mathbf{w} = \arg \min_{\mathbf{w}} \sum_{n=1}^N d(y_n, f(\mathbf{x}_n))$$

$$g_{\text{logistic}}(z) = \frac{e^z}{1 + e^z} = \frac{1}{1 + e^{-z}}$$

$g_{\text{logistic}}(z)$



Interpretation of $f(\mathbf{x})$: The probability that the observation belongs to class 1

Generalized linear model

- **Cost function:** Choose one

$$d(y, f(\mathbf{x}))$$

- **Model:** Linear + non-linear link

$$f(\mathbf{x}) = g_{\text{link}}(\mathbf{x}^\top \mathbf{w})$$

- **Parameters:** Optimize using numerical optimization methods

$$\mathbf{w} = \arg \min_{\mathbf{w}} \sum_{n=1}^N d(y_n, f(\mathbf{x}_n))$$

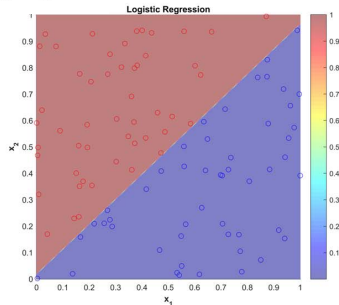
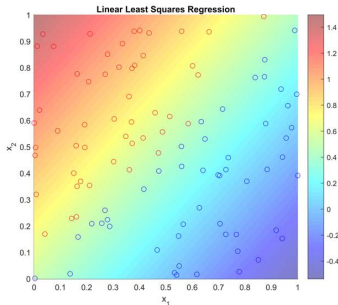
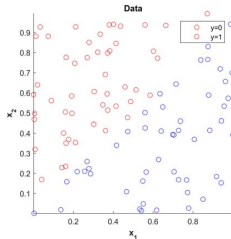


Matlab: `glmfit`
 Python: `sklearn.linear_model`
 R: `glm`

$$g_{\text{identity}}(z) = z$$

$$g_{\text{logistic}}(z) = \frac{e^z}{1 + e^z} = \frac{1}{1 + e^{-z}}$$

Linear vs. Logistic Regression for a classification problem



$$f(x) = 0.5270 - 1.0639x_1 + 0.9684x_2$$

$$f(x) = 1 / [1 + \exp(-(-41.6 - 2295x_1 + 2422x_2))]$$