

მონაცემთა ანალიტიკა Python

ლექცია 14: კლასიფიკაცია. ერთ-ცვლადიანი ლოჯისტიკური რეგრესია. მრავალ-ცვლადიანი ლოჯისტიკური რეგრესია

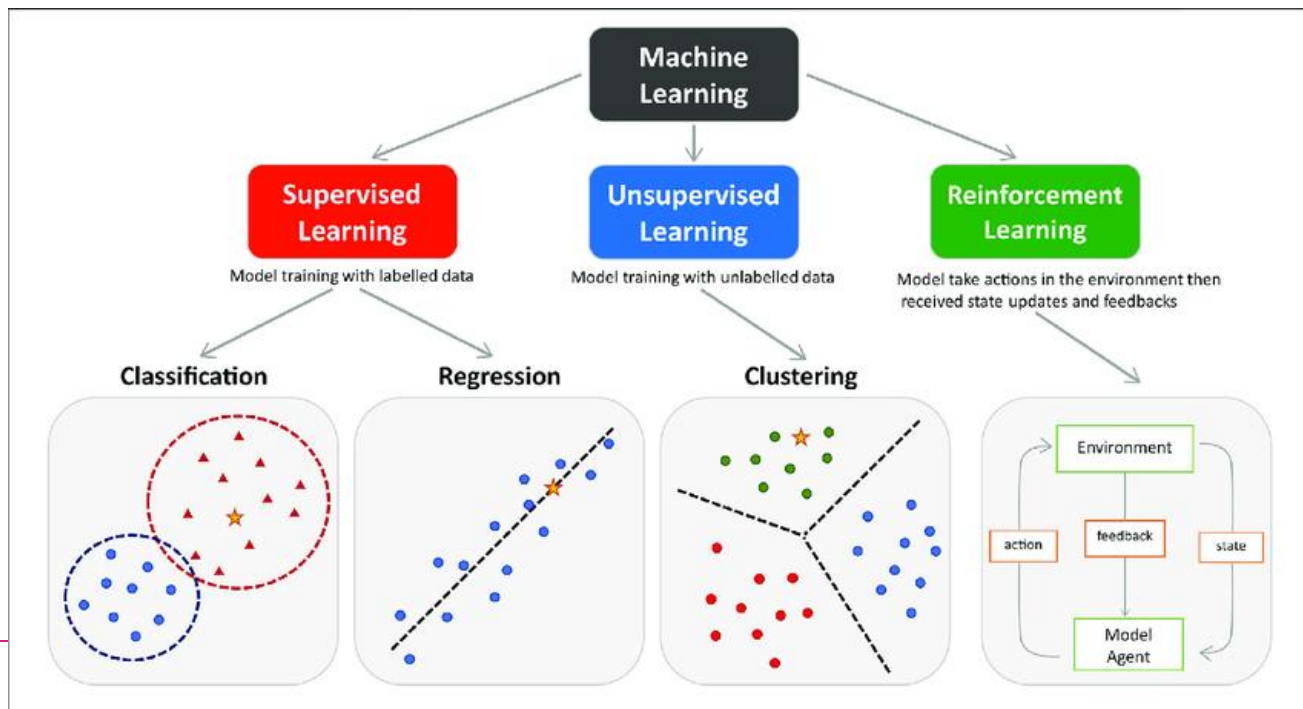
ლექცია 15: კლასიფიკაცია. გადაწყვეტილების ხე

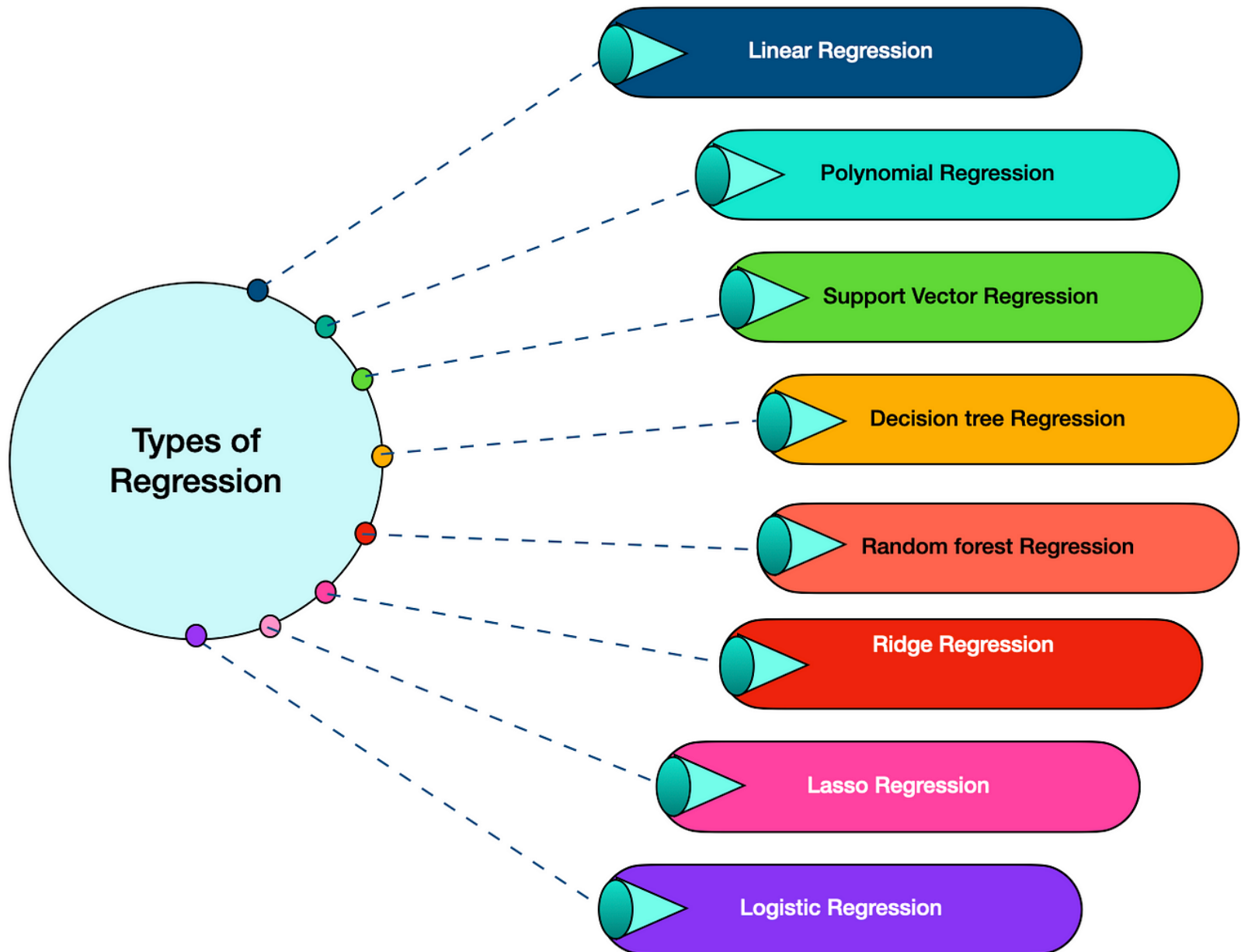
ლიკა სვანაძე
lika.svanadze@btu.edu.ge

Machine learning Models

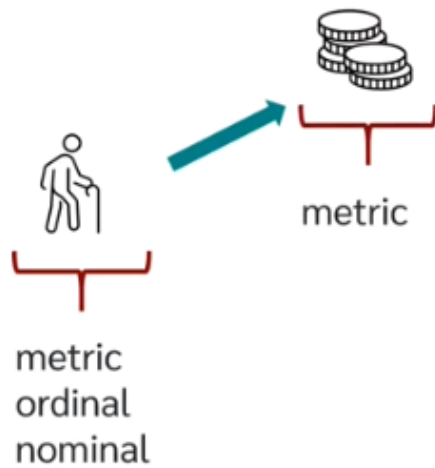
Based on the tasks performed and the nature of the output, you can **classify** machine learning models into three types:

1. **Regression:** where the output variable to be predicted is a continuous variable
2. **Classification:** where the output variable to be predicted is a categorical variable
3. **Clustering:** where there is no pre-defined notion of a label allocated to the groups/clusters formed.

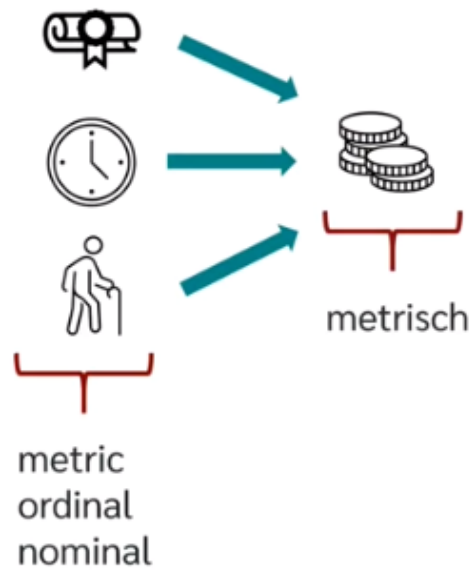




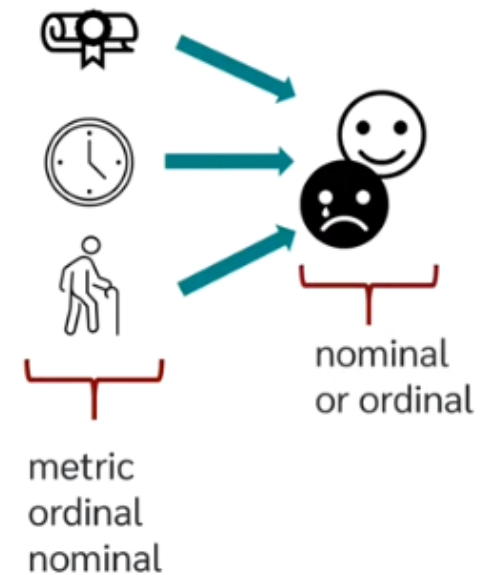
Simple linear Regression



Multiple linear Regression

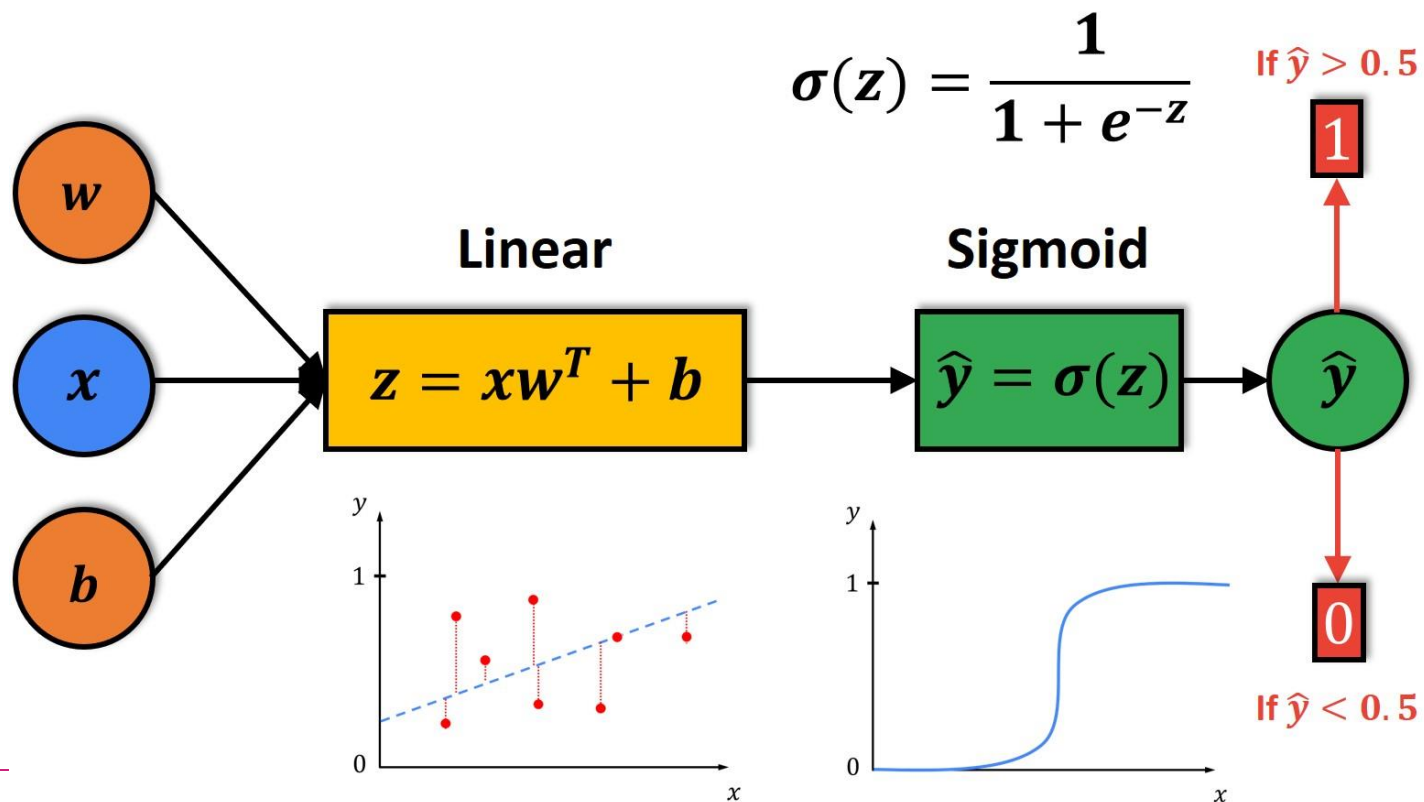


Logistic Regression



Logistic Regression

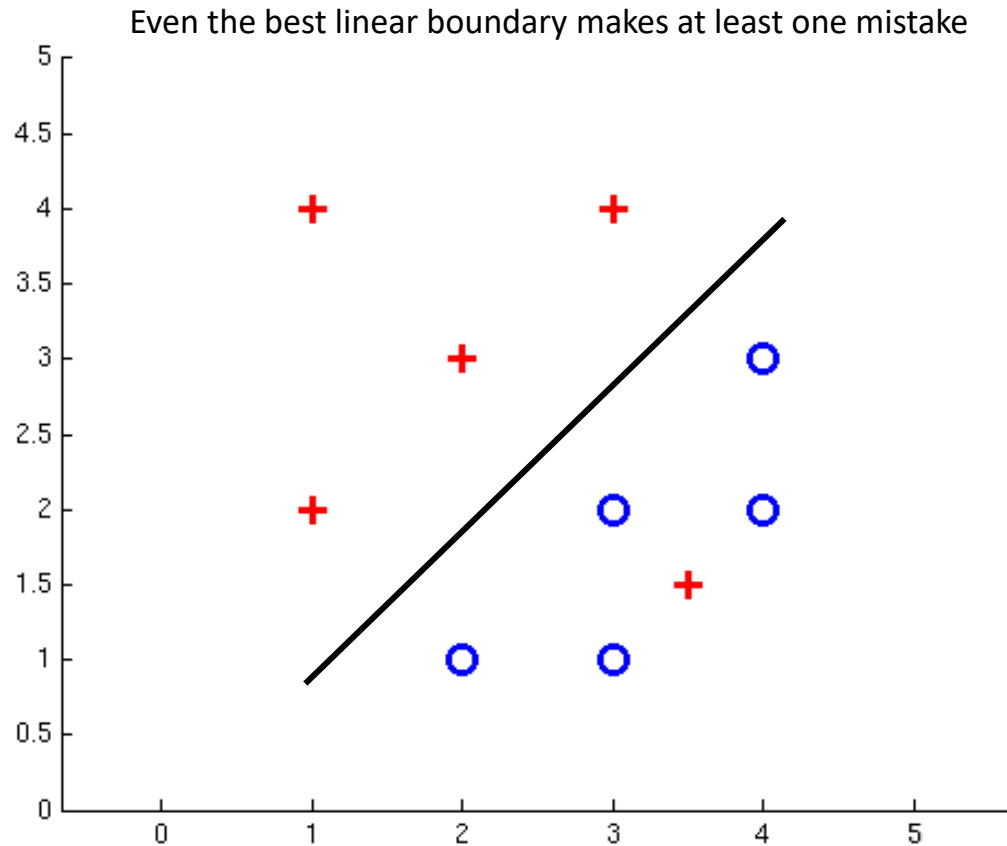
- Unlike linear regression which outputs continuous number values, **logistic regression** uses the logistic sigmoid function to transform its output to return a probability value which can then be mapped to two or more discrete classes.



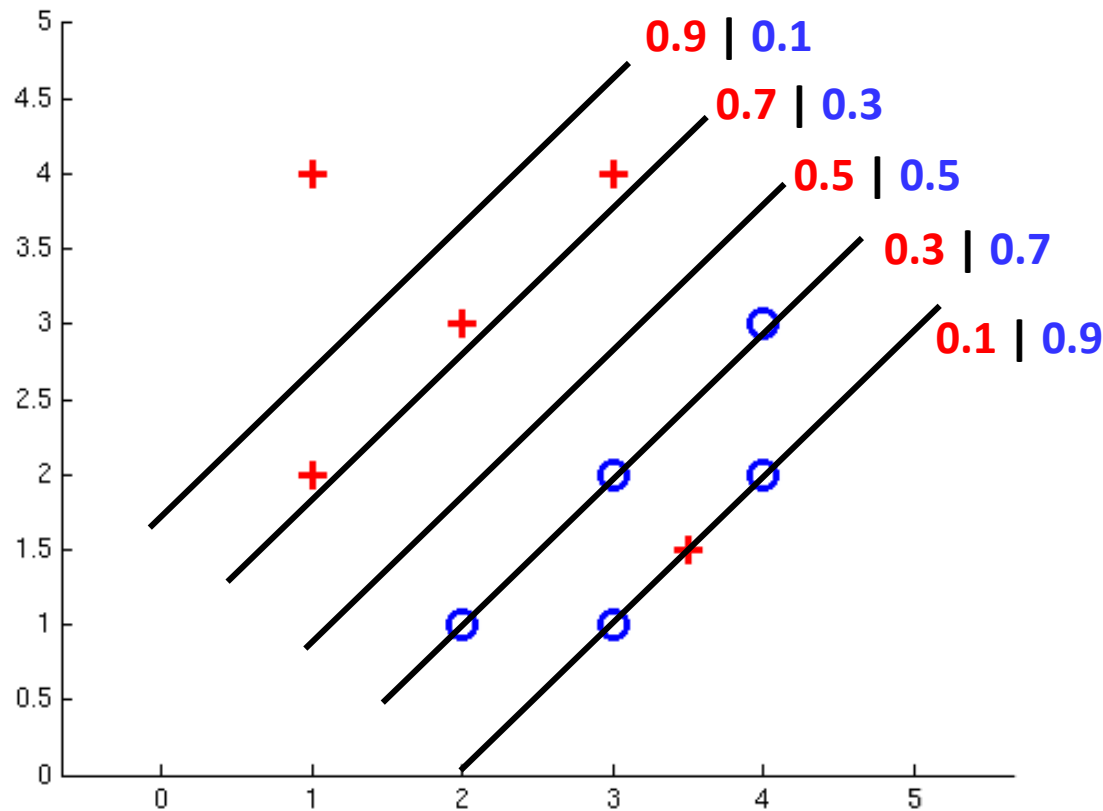
Types of Logistic Regression

1. Binary (true/false, yes/no)
2. Multi-class (sheep, cats, dogs)
3. Ordinal (Job satisfaction level — dissatisfied, satisfied, highly satisfied)

Non-Separable Case: Deterministic Decision



Non-Separable Case: Probabilistic Decision



How to get probabilistic decisions?

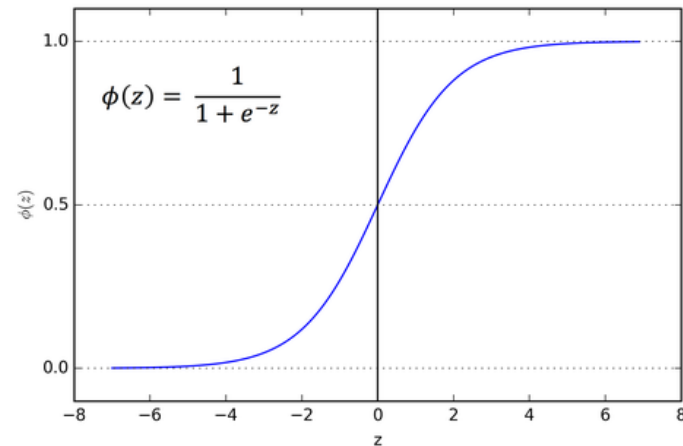
Perceptron scoring: $z = w \cdot f(x)$

If $z = w \cdot f(x)$ very positive \rightarrow want probability going to 1

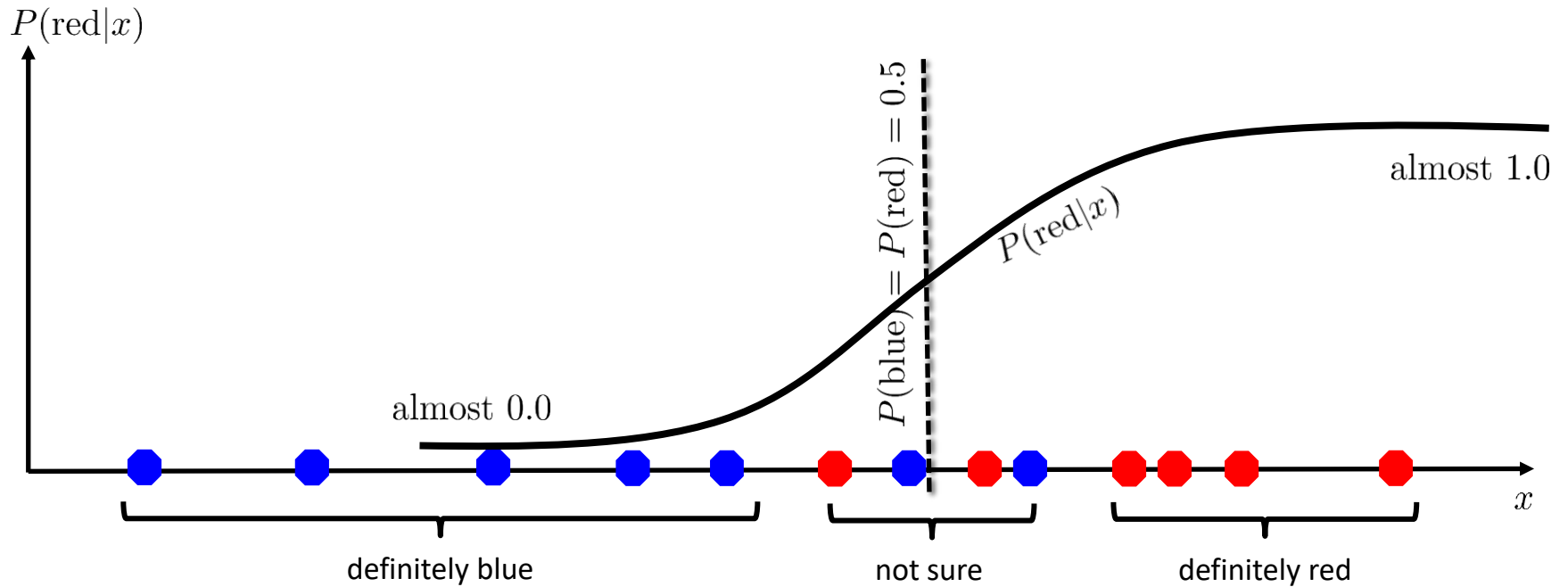
If $z = w \cdot f(x)$ very negative \rightarrow want probability going to 0

Sigmoid function

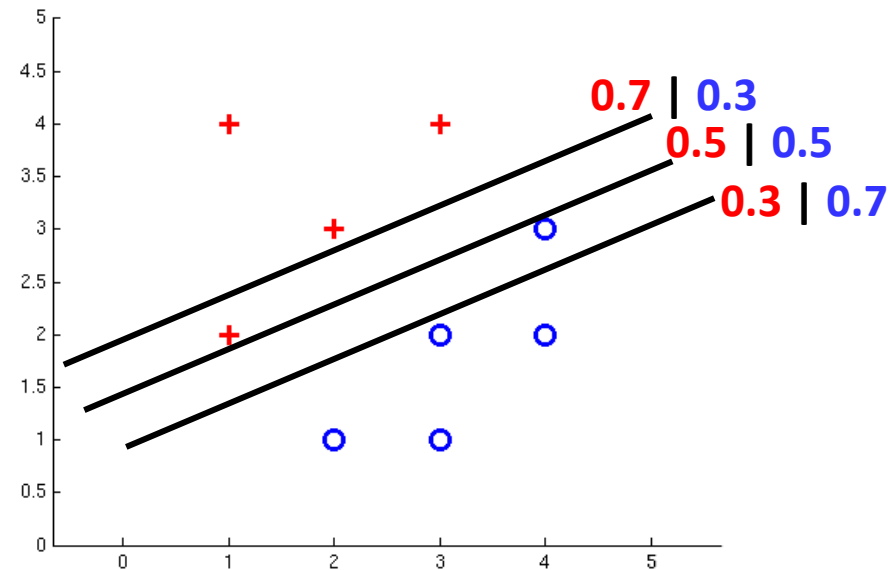
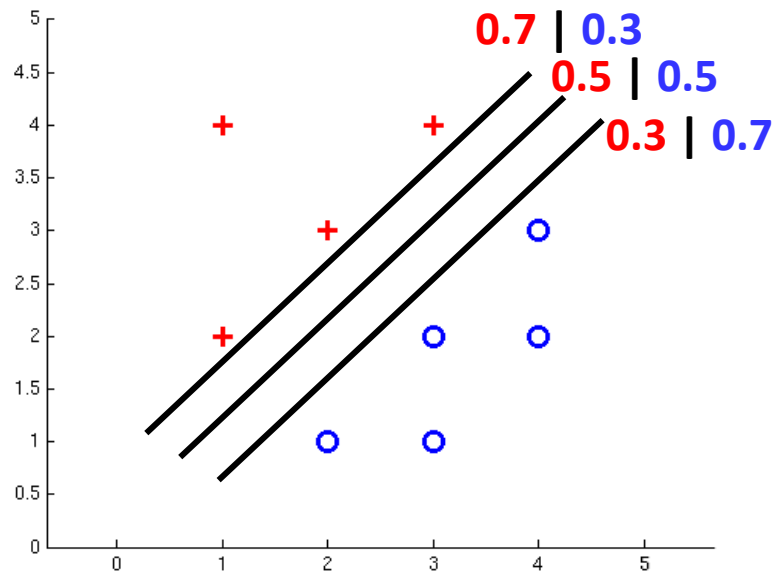
$$\phi(z) = \frac{1}{1 + e^{-z}}$$



A 1D Example



Separable Case: Probabilistic Decision – Clear Preference



Best w ?

Maximum likelihood estimation:

$$\max_w ll(w) = \max_w \sum_i \log P(y^{(i)} | x^{(i)}; w)$$

with:

$$P(y^{(i)} = +1 | x^{(i)}; w) = \frac{1}{1 + e^{-w \cdot f(x^{(i)})}}$$

$$P(y^{(i)} = -1 | x^{(i)}; w) = 1 - \frac{1}{1 + e^{-w \cdot f(x^{(i)})}}$$

= Logistic Regression

მონაცემთა ანალიტიკა Python

ლექცია 15: კლასიფიკაცია. გადაწყვეტილების ხე. მრავალი
გადაწყვეტილების ხე.

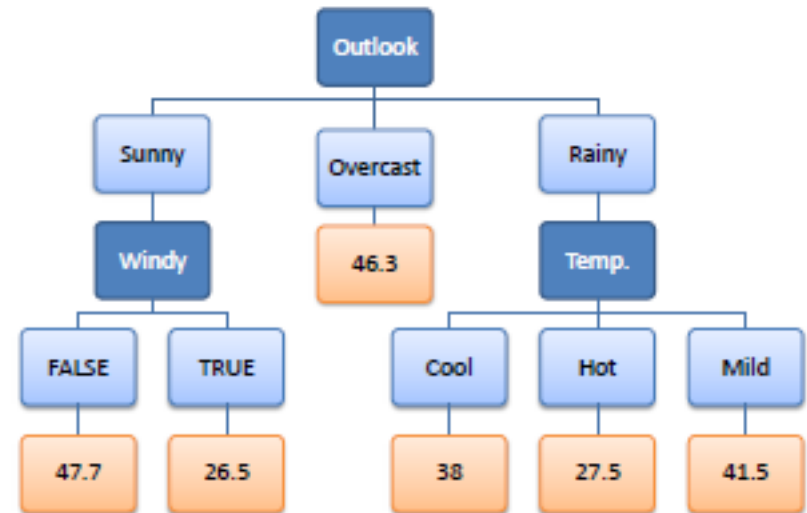
ლიკა სვანაძე
lika.svanadze@btu.edu.ge

Decision Trees

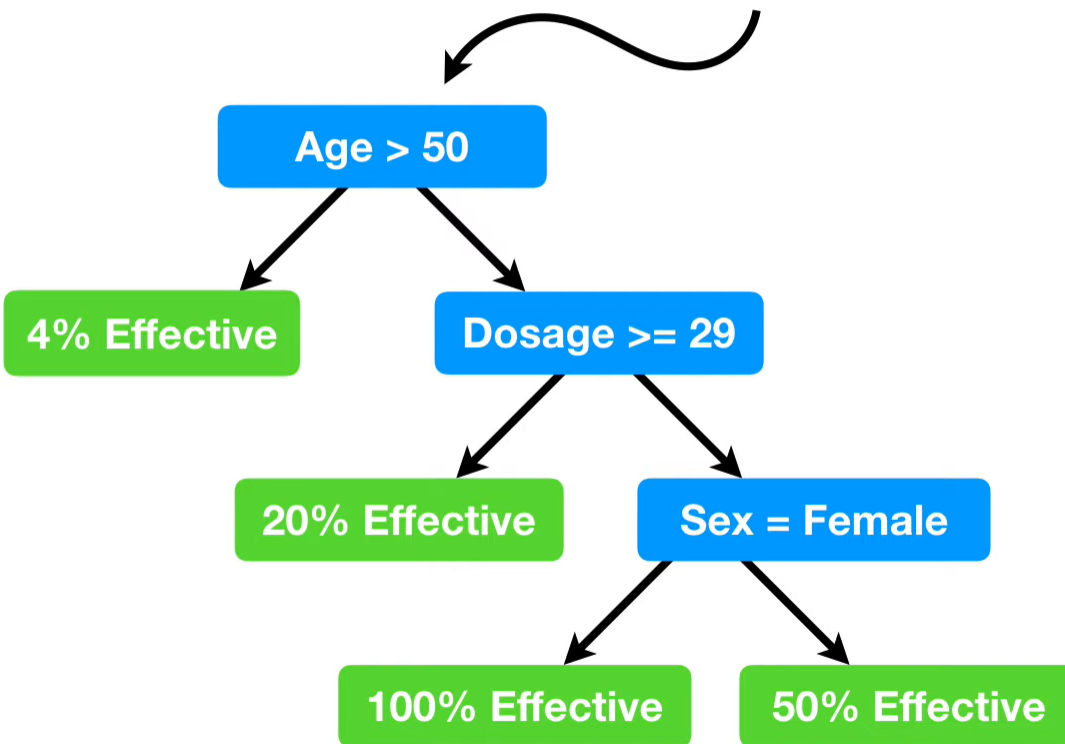


Decision Tree - Regression

Predictors				Target
Outlook	Temp.	Humidity	Windy	Hours Played
Rainy	Hot	High	False	26
Rainy	Hot	High	True	30
Overcast	Hot	High	False	48
Sunny	Mild	High	False	46
Sunny	Cool	Normal	False	62
Sunny	Cool	Normal	True	23
Overcast	Cool	Normal	True	43
Rainy	Mild	High	False	36
Rainy	Cool	Normal	False	38
Sunny	Mild	Normal	False	48
Rainy	Mild	Normal	True	48
Overcast	Mild	High	True	62
Overcast	Hot	Normal	False	44
Sunny	Mild	High	True	30



In contrast, a **Regression Tree** easily accommodates the additional predictors.



Dosage	Age	Sex	Etc.	Drug Effect.
10	25	Female	...	98
20	73	Male	...	0
35	54	Female	...	100
5	12	Male	...	44
etc...	etc...	etc...	etc...	etc...

Decision Trees



Reminder: Features

Features, aka attributes

- Sometimes: $\text{TYPE}=\text{French}$
- Sometimes: $f_{\text{TYPE}=\text{French}}(x) = 1$

Example	Attributes										Target
	<i>Alt</i>	<i>Bar</i>	<i>Fri</i>	<i>Hun</i>	<i>Pat</i>	<i>Price</i>	<i>Rain</i>	<i>Res</i>	<i>Type</i>	<i>Est</i>	<i>WillWait</i>
X_1	<i>T</i>	<i>F</i>	<i>F</i>	<i>T</i>	<i>Some</i>	<i>\$\$\$</i>	<i>F</i>	<i>T</i>	<i>French</i>	<i>0-10</i>	<i>T</i>
X_2	<i>T</i>	<i>F</i>	<i>F</i>	<i>T</i>	<i>Full</i>	<i>\$</i>	<i>F</i>	<i>F</i>	<i>Thai</i>	<i>30-60</i>	<i>F</i>
X_3	<i>F</i>	<i>T</i>	<i>F</i>	<i>F</i>	<i>Some</i>	<i>\$</i>	<i>F</i>	<i>F</i>	<i>Burger</i>	<i>0-10</i>	<i>T</i>
X_4	<i>T</i>	<i>F</i>	<i>T</i>	<i>T</i>	<i>Full</i>	<i>\$</i>	<i>F</i>	<i>F</i>	<i>Thai</i>	<i>10-30</i>	<i>T</i>
X_5	<i>T</i>	<i>F</i>	<i>T</i>	<i>F</i>	<i>Full</i>	<i>\$\$\$</i>	<i>F</i>	<i>T</i>	<i>French</i>	<i>>60</i>	<i>F</i>
X_6	<i>F</i>	<i>T</i>	<i>F</i>	<i>T</i>	<i>Some</i>	<i>\$\$</i>	<i>T</i>	<i>T</i>	<i>Italian</i>	<i>0-10</i>	<i>T</i>
X_7	<i>F</i>	<i>T</i>	<i>F</i>	<i>F</i>	<i>None</i>	<i>\$</i>	<i>T</i>	<i>F</i>	<i>Burger</i>	<i>0-10</i>	<i>F</i>
X_8	<i>F</i>	<i>F</i>	<i>F</i>	<i>T</i>	<i>Some</i>	<i>\$\$</i>	<i>T</i>	<i>T</i>	<i>Thai</i>	<i>0-10</i>	<i>T</i>
X_9	<i>F</i>	<i>T</i>	<i>T</i>	<i>F</i>	<i>Full</i>	<i>\$</i>	<i>T</i>	<i>F</i>	<i>Burger</i>	<i>>60</i>	<i>F</i>
X_{10}	<i>T</i>	<i>T</i>	<i>T</i>	<i>T</i>	<i>Full</i>	<i>\$\$\$</i>	<i>F</i>	<i>T</i>	<i>Italian</i>	<i>10-30</i>	<i>F</i>
X_{11}	<i>F</i>	<i>F</i>	<i>F</i>	<i>F</i>	<i>None</i>	<i>\$</i>	<i>F</i>	<i>F</i>	<i>Thai</i>	<i>0-10</i>	<i>F</i>
X_{12}	<i>T</i>	<i>T</i>	<i>T</i>	<i>T</i>	<i>Full</i>	<i>\$</i>	<i>F</i>	<i>F</i>	<i>Burger</i>	<i>30-60</i>	<i>T</i>

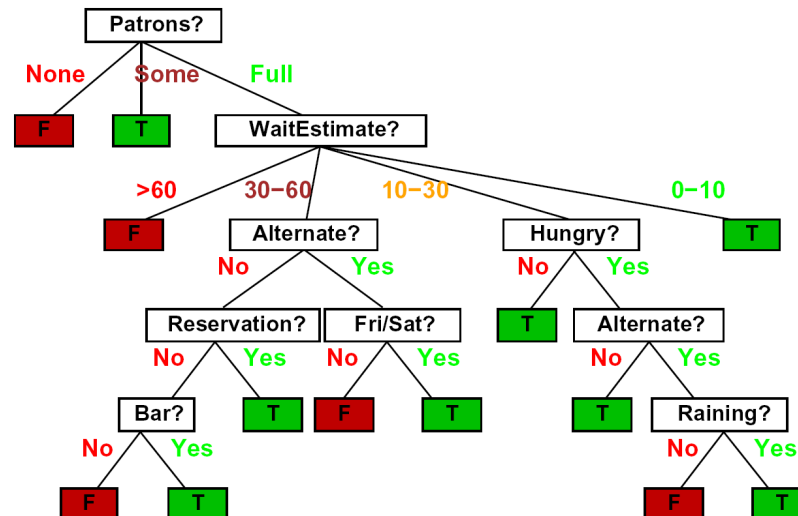
Decision Trees

Compact representation of a function:

- Truth table
- Conditional probability table
- Regression values

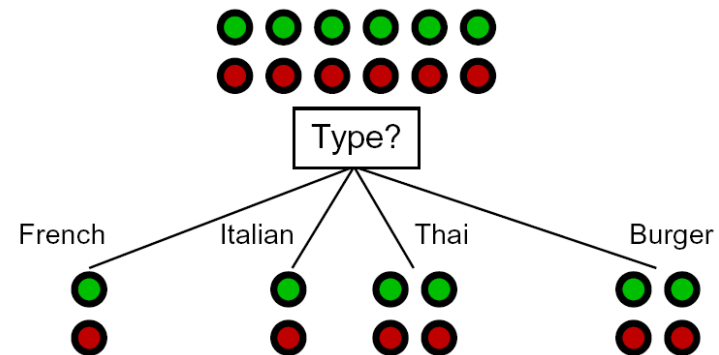
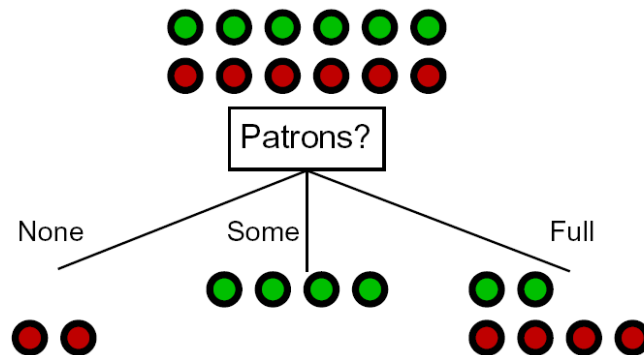
True function

- Realizable: in H



Choosing an Attribute

Idea: a good attribute splits the examples into subsets that are (ideally) “all positive” or “all negative”



So: we need a measure of how “good” a split is, even if the results aren’t perfectly separated out

Entropy and Information

Information answers questions

- The more uncertain about the answer initially, the more information in the answer
- Scale: bits
 - Answer to Boolean question with prior $\langle 1/2, 1/2 \rangle$?
 - Answer to 4-way question with prior $\langle 1/4, 1/4, 1/4, 1/4 \rangle$?
 - Answer to 4-way question with prior $\langle 0, 0, 0, 1 \rangle$?
 - Answer to 3-way question with prior $\langle 1/2, 1/4, 1/4 \rangle$?

A probability p is typical of:

- A uniform distribution of size $1/p$
- A code of length $\log 1/p$

Entropy

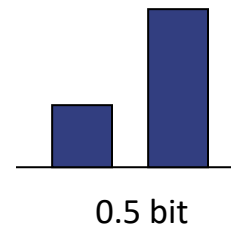
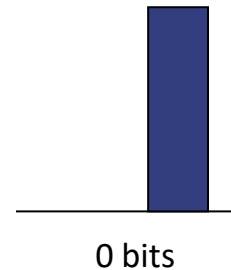
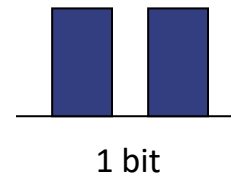
General answer: if prior is $\langle p_1, \dots, p_n \rangle$:

- Information is the expected code length

$$\begin{aligned} H(\langle p_1, \dots, p_n \rangle) &= E_p \log_2 1/p_i \\ &= \sum_{i=1}^n -p_i \log_2 p_i \end{aligned}$$

Also called the **entropy** of the distribution

- More uniform = higher entropy
- More values = higher entropy
- More peaked = lower entropy
- Rare values almost “don’t count”

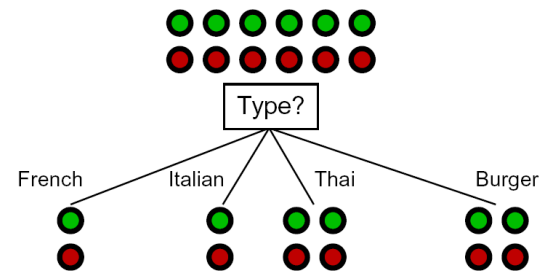
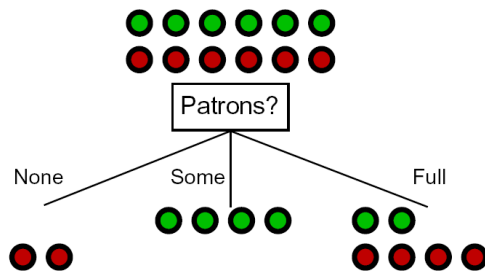


Information Gain

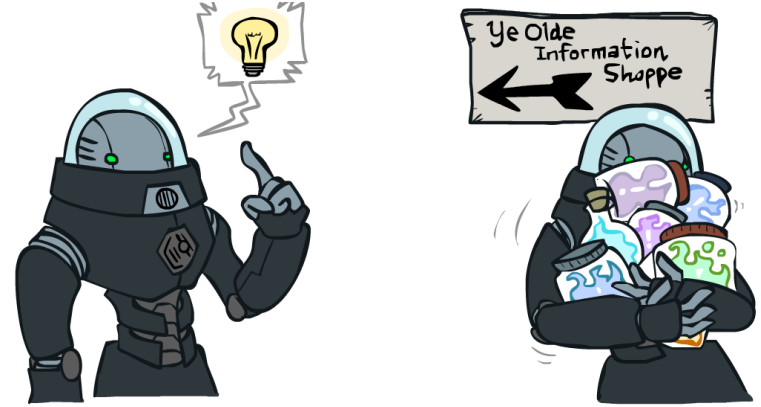
Back to decision trees!

For each split, compare entropy before and after

- Difference is the **information gain**
- Problem: there's more than one distribution after split!



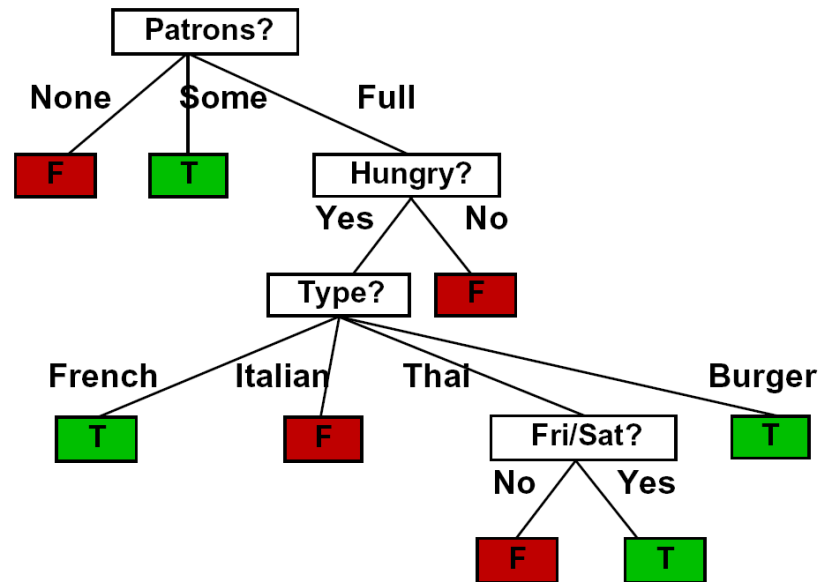
- Solution: use **expected entropy**, weighted by the number of examples



BTU | გიგანტურა და
ხანგრძლივობა
ენიკანობა

Example: Learned Tree

Decision tree learned from these 12 examples:



Substantially simpler than “true” tree

- A more complex hypothesis isn't justified by data

Also: it's reasonable, but wrong

Python implementation (See python files)