# Classification

## with Decision Trees

## MODALG SoSe18

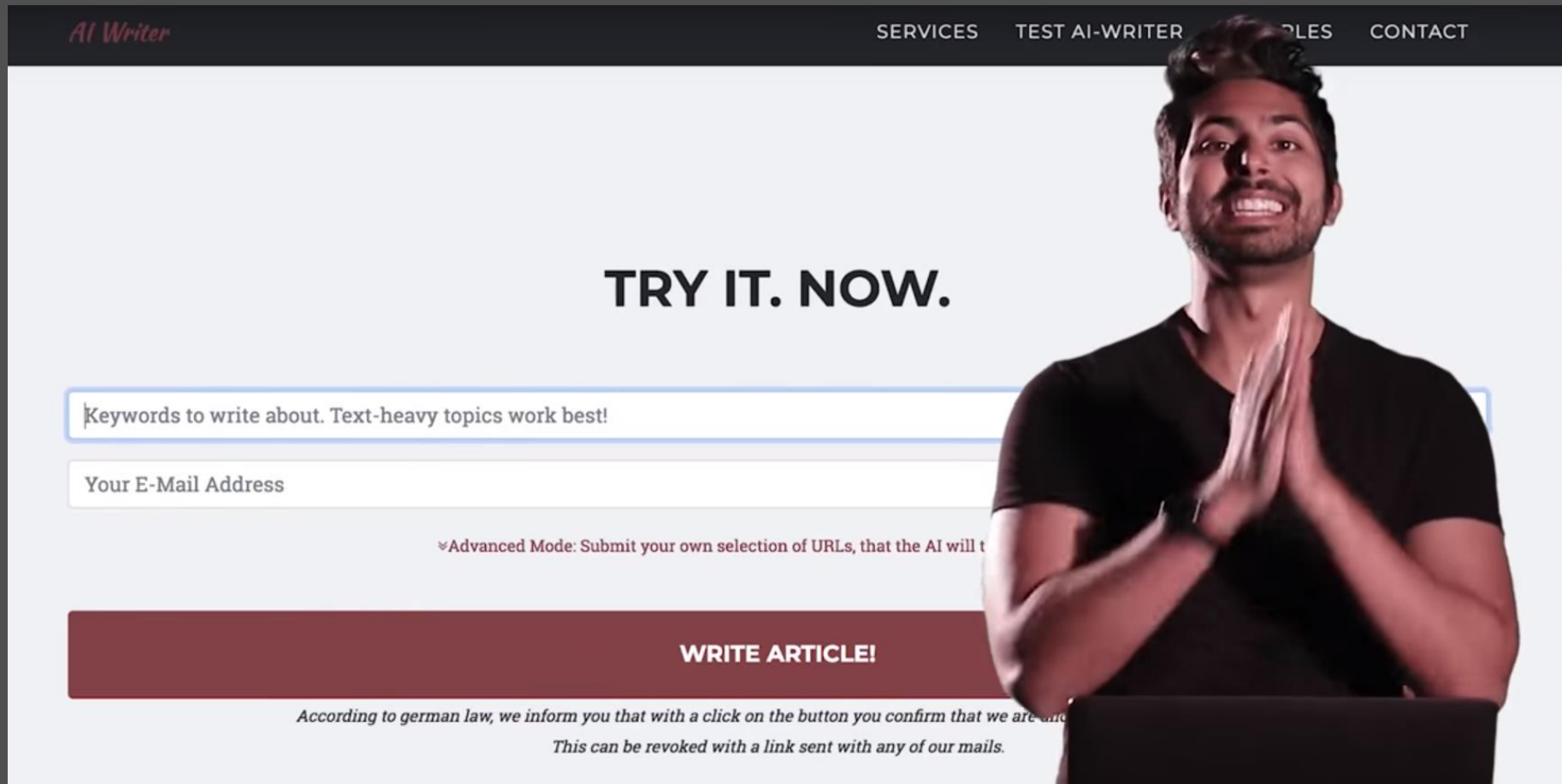# Facts & News

# OpenAI Salaries (Source NYT):

- **Ilya Sutskever: 1.9 MM $ (in 2016)**
- **Ian Goodfellow: 800 M $ (in 2016)**
- **Pieter Abbeel: 425 M $ (in 2016)**

# AI in Marketing

# Start of Google I/O 2018

# Repitition

# Artificial Intelligence

**Artificial Intelligence**

**Machine Learning**

Artificial Intelligence

Machine Learning

Neural Nets

"A computer program is said to learn from **experience E with respect to some class of tasks T and performance measure P**, if its performance at tasks in T, as measured by P, improves with experience E." (Tom Mitchell)

# Human Perspective:

- Task: Learn a text by heart!

- Performance: How many words are wrong?

- Experience: Learning Progress

Answer

Pupil → Teacher

Correction

Prediction

Predictive Model → Ground Truth Data

Error

# Classification

# 5 Questions

What is the aim of classification?

What is a class?

Natural classification?

Artificial classification?

How we can do it?

# What is the aim of classification?

- **Differentiation**

# What is a class?

- Result of grouping things
- Compare with different types of personalities
- Data with the specific label

# Natural classification?

- **You do it daily!!!**
- **Learn something about the world**
- **Apply knowledge**

# Artificial classification?

- Let an algorithm learn something about the world
- Apply knowledge

# How we can do it?

- Data (knowledge about the world)
- Algorithm (something that can learn)

# In Detail

# Example Class Survived

# Can we predict if someone would survive?
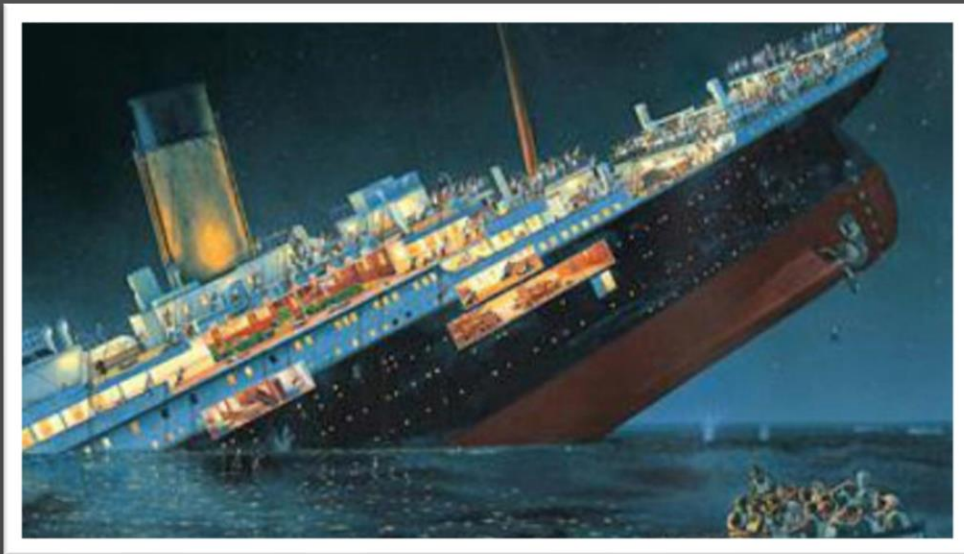
- Look at the features!
  - Male or female
  - Number of class
  - Age of the person

angelina_jolie.21 angelina_jolie.22 angelina_jolie.23 angelina_jolie.24 angelina_jolie.25 angelina_jolie.26 angelina_jolie.27

angelina_jolie.28 angelina_jolie.29 angelina_jolie.30 angelina_jolie.31 angelina_jolie.32 angelina_jolie.3

Elementtyp: JPG-Datei
Bewertung: Nicht bewert
Abmessungen: 1200 x 68
Größe: 81,8 KB

angelina_jolie.35 angelina_jolie.36 angelina_jolie.37 angelina_jolie.38 angelina_jolie.39 angelina_jolie.40 angelina_jolie.41

angelina_jolie.42 angelina_jolie.44 angelina_jolie.45 angelina_jolie.46 angelina_jolie.47 angelina_jolie.48 angelina_jolie.49

angelina_jolie.50 angelina_jolie.51 angelina_jolie.52 angelina_jolie.53 angelina_jolie.54 angelina_jolie.55 angelina_jolie.56

# Algorithm

## Algorithm

**Decision Tree**

## Algorithm

**Random Forest**

**Decision Tree**

**Gradient Boosting**

**Algorithm**

**Random Forest**

Naive Bayes

Decision Tree

Gradient Boosting

Algorithm

Random Forest

Naive Bayes

Decision Tree

Gradient Boosting

Gaussian Mixture Model

# Algorithm

Random Forest

**Naive Bayes**

**Decision Tree**

**Gradient Boosting**

**Algorithm**

**Gaussian Mixture Model**

**Random Forest**

**Support Vector Machine**

Naive Bayes

Decision Tree

Gradient Boosting

Gaussian Mixture Model

**Algorithm**

Random Forest

K-Nearest-Neighbor

Support Vector Machine

Naive Bayes          Decision Tree

Gradient Boosting          Gaussian Mixture Model

Neural Network

**Algorithm**

Random Forest          Support Vector Machine

K-Nearest-Neighbor

# Algorithm = Model

# Decision Tree

- **Should I eat sweets?**
- **Possible answers: Yes/No**
- **Things to consider:**
  - Weight
  - Condition of teeth
  - Diabetes
  - Time
  - Situation

# Questions and Splits

- **Given a dataset**
- **Assume we have some questions concerning the data**
- **How to decide where to split?**

- **Find the feature that best splits the target class into the purest possible children nodes**

- **If target class consist male and female**

- **Nodes that don't contain a mix of both male and female, rather pure nodes with only one class**

# Information Gain

- **Information Gain is a measurement**

- **How much information do we gain by doing a split at particular feature?**

- **Compare to measure of quality when you do a split at specific question**

# Information Gain

$$\text{information gain} = \text{entropy (parent)} - \left[\begin{array}{c}\text{weighted} \\ \text{average}\end{array}\right] \text{entropy (children)}$$

decision tree algorithm : maximize information gai

# What is entropy?

- Another measurement
- Measures **impurity** of data

$$H(x) = -\sum_{i=1}^{n} p(x_i) * log_b(p(x_i))$$

$H$:    **Entropy**

$x$:    **Whole column of target feature (e.g. Eat sweets?)**

$x_i \in \{0,1\}$

**e.g.**  $x = (0,1,0,0,1,0,1,0,1)^T$

$$H(x) = -\sum_{i=1}^{n} p(x_i) * log_b(p(x_i))$$

$p(x_i)$:     **How many times is** $x_i = 0$ **? (relative value)**
$p(x_i)$:     **How many times is** $x_i = 1$ **? (relative value)**
$x_i$:     **Answer of eating sweets**

$$H(x) = -\sum_{i=1}^{n} p(x_i) * log_b(p(x_i))$$

$b$: **number of classes (binary:** $b = 2$ **)**

$n$: **length of vector** $x$

$$H(x) = -\sum_{i=1}^{n} p(x_i) * log_b(p(x_i))$$

# What else?

- **Recursive algorithm (Iterative Dichotomiser - ID3)**

- **Assume binary problem:**

**build_tree(data, questions){**

**}**

- **Recursive algorithm (Iterative Dichotomiser - ID3)**

- **Assume binary problem:**

```
build_tree(data, questions){
    foreach(question) gain, question = information_gain(data, questions)
}
```

- **Recursive algorithm (Iterative Dichotomiser - ID3)**

- **Assume binary problem:**

```
build_tree(data, questions){
    foreach(question) gain, question = information_gain(data, questions)
    take question with argmax gain
}
```

- ## Recursive algorithm (Iterative Dichotomiser - ID3)

- ## Assume binary problem:

```
build_tree(data, questions){
    foreach(question) gain, question = information_gain(data, questions)
    take question with argmax gain
    if(gain == 0):
        return leaf_node(data)
}
```

- ## Recursive algorithm (Iterative Dichotomiser - ID3)

- ## Assume binary problem:

```
build_tree(data, questions){
    foreach(question) gain, question = information_gain(data, questions)
    take question with argmax gain
    if(gain == 0):
        return leaf_node(data)
    true_data, false_data = split(data, question)
}
```

- **Recursive algorithm (Iterative Dichotomiser - ID3)**

- **Assume binary problem:**

```
build_tree(data, questions){
    foreach(question) gain, question = information_gain(data, questions)
    take question with argmax gain
    if(gain == 0):
        return leaf_node(data)
    true_data, false_data = split(data, question)
    build_tree(true_data, questions)
    build_tree(false_data, questions)
}
```

- ## Recursive algorithm (Iterative Dichotomiser - ID3)

- ## Assume binary problem:

```
build_tree(data, questions){
    foreach(question) gain, question = information_gain(data, questions)
    take question with argmax gain
    if(gain == 0):
        return leaf_node(data)
    true_data, false_data = split(data, question)
    build_tree(true_data, questions)
    build_tree(false_data, questions)
    return decision_node(question, data)
}
```
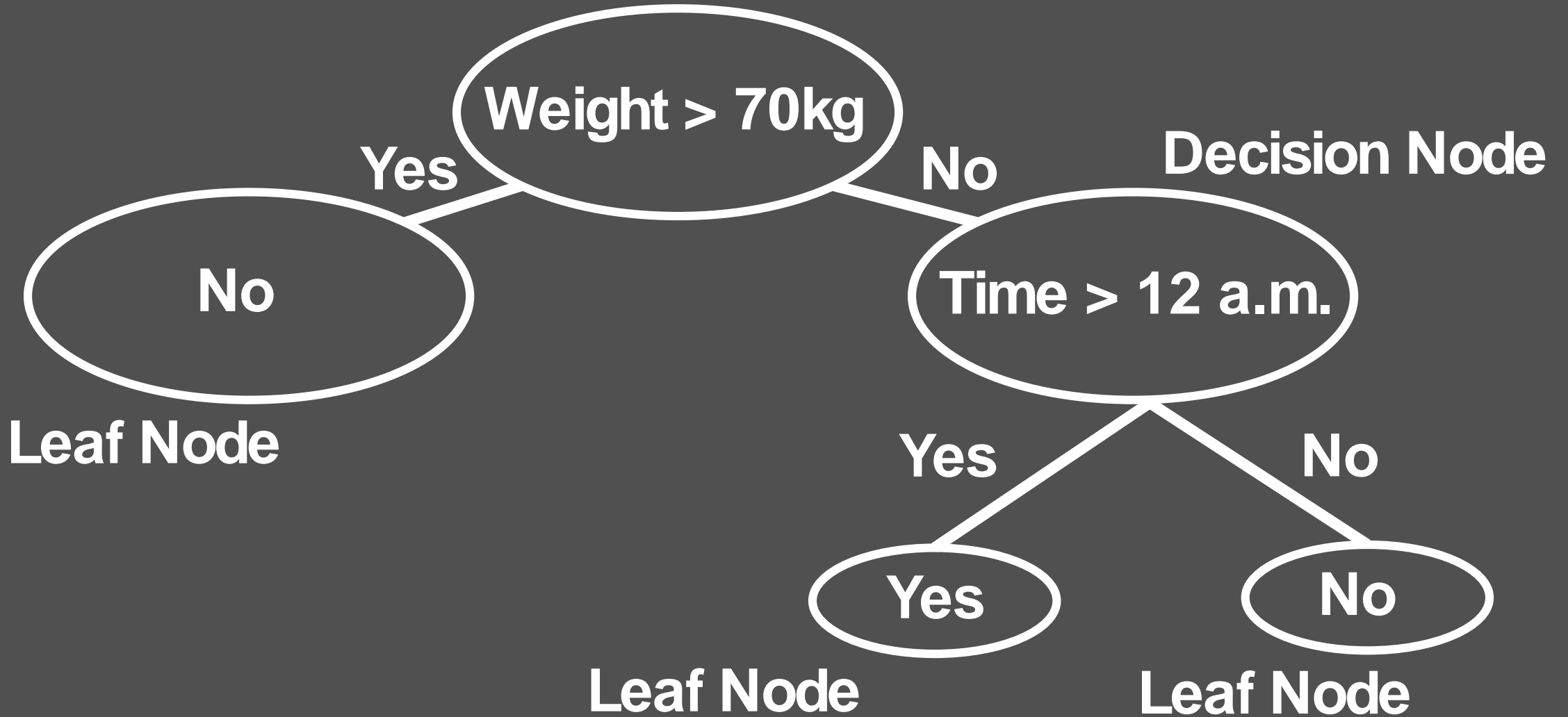
# There is more!

- **Multiclass-Problem (Eat sweets? Yes/No/Maybe)**
- **Pruning of the tree**
- **Cross Validation**
- **Bias**
- **Metrics – Important: Accuracy**

# Accuracy = #Correct Classified / #Examples
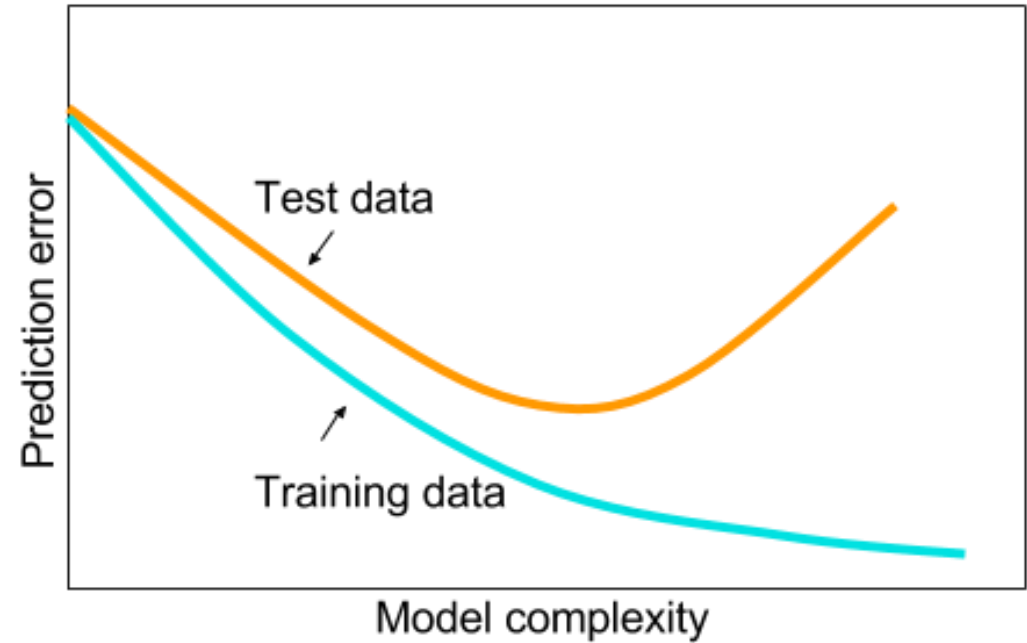
# What is the aim of classification?

- **Differentiation**
- **You want to be a master**
- **Generalise!**

# Problem of generalisation:

- The world is full of **bias**

- You may a master concerning the training set, but what happens if I give you a another test set?

- Accuracy can decrease!

# Overfitting

# Conclusion:

- **You can compare classification to grouping**
- **We can use algorithms (models)**
- **A decision tree is a human readable model**
- **Measure the performance with classification accuracy**

# Your Turn!

# github.com/mati3230/modalg181