# Can Machine Learning solve my problem?

*9th BCD ISHEP Cargèse School - part I*

*28 March 2024 - Cargèse, France*

## Emille E. O. Ishida

*Laboratoire de Physique de Clermont - Université Clermont-Auvergne*
*Clermont Ferrand, France*

*What **impressive things** machine learning can and/or will be able to do?*

# Can Machine Learning solve my problem?

*I do not know*

*9th BCD ISHEP Cargèse School  - part I*
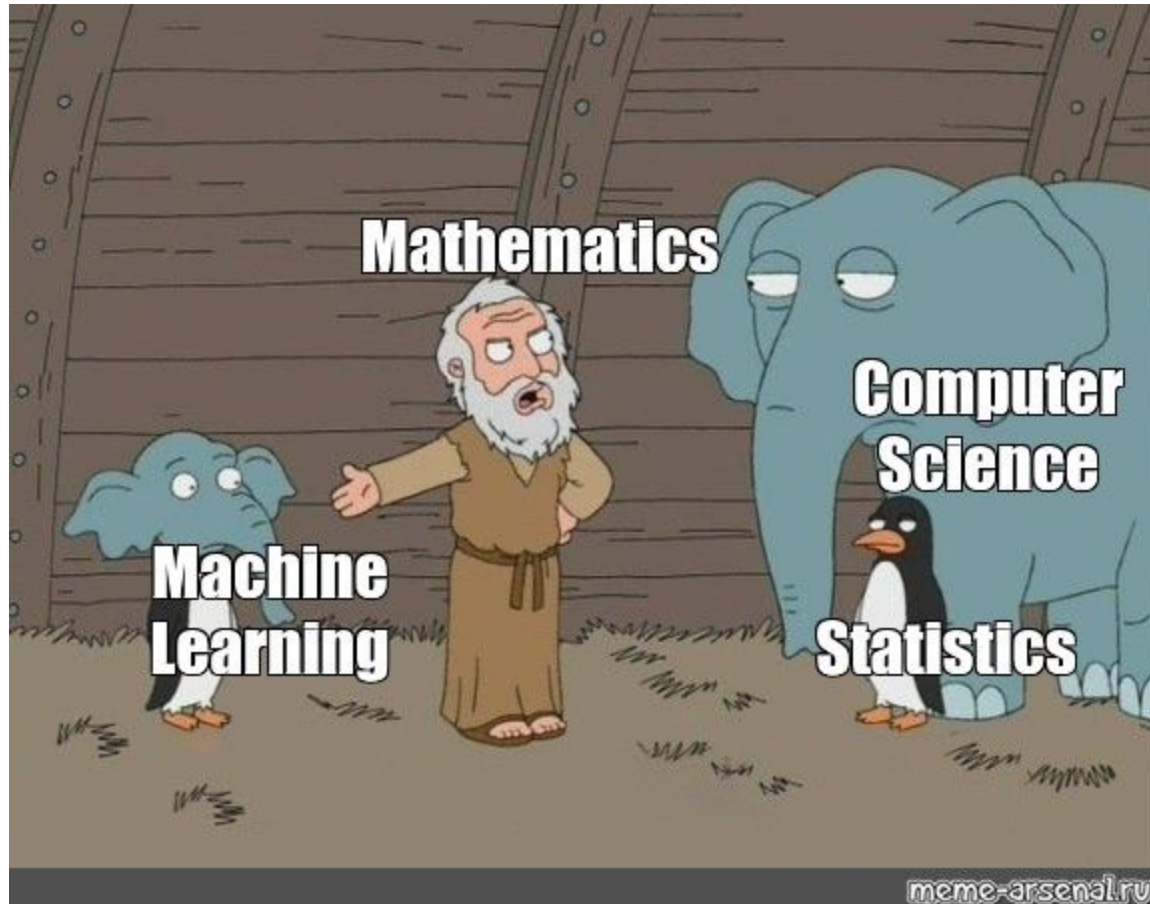*28 March 2024 - Cargèse, France*

## Emille E. O. Ishida

*Laboratoire de Physique de Clermont - Université Clermont-Auvergne*
*Clermont Ferrand, France*

# What is learning ?

*"A relatively permanent change in behaviour due to past experiences."*

D. Coon, *Introduction to psychology: exploration and application* (1983)

# Supervised Learning

*Learn by example*

**Learn**

**Training sample**

# Features
# +
# Labels

⟹

**Apply**

**Target sample**

Features

# Question:

List 2 animals that you believe are capable of learning.



Join at menti.com with code: 7907 6385

# Rat bait shyness - I

*Rzóska, J. (1953). Bait shyness, a study in rat behavior. British Journal of Animal Behaviour, 1, 128–135*

# Rat bait shyness - II

*Rzóska, J. (1953). Bait shyness, a study in rat behavior. British Journal of Animal Behaviour, 1, 128–135*

# Question:

- Do you believe the rat will learn the correlation between bad food ⇒ shock and/or sound ⇒ nausea?

Join at menti.com with code: 7907 6385

# Question:

- What aspect of the rat learning model prevents it from understanding the input ⇒ output correlation?

# Pigeon superstition

*Skinner, B. F. "'Superstition' in the Pigeon", Journal of Experimental Psychology#38, 1947*
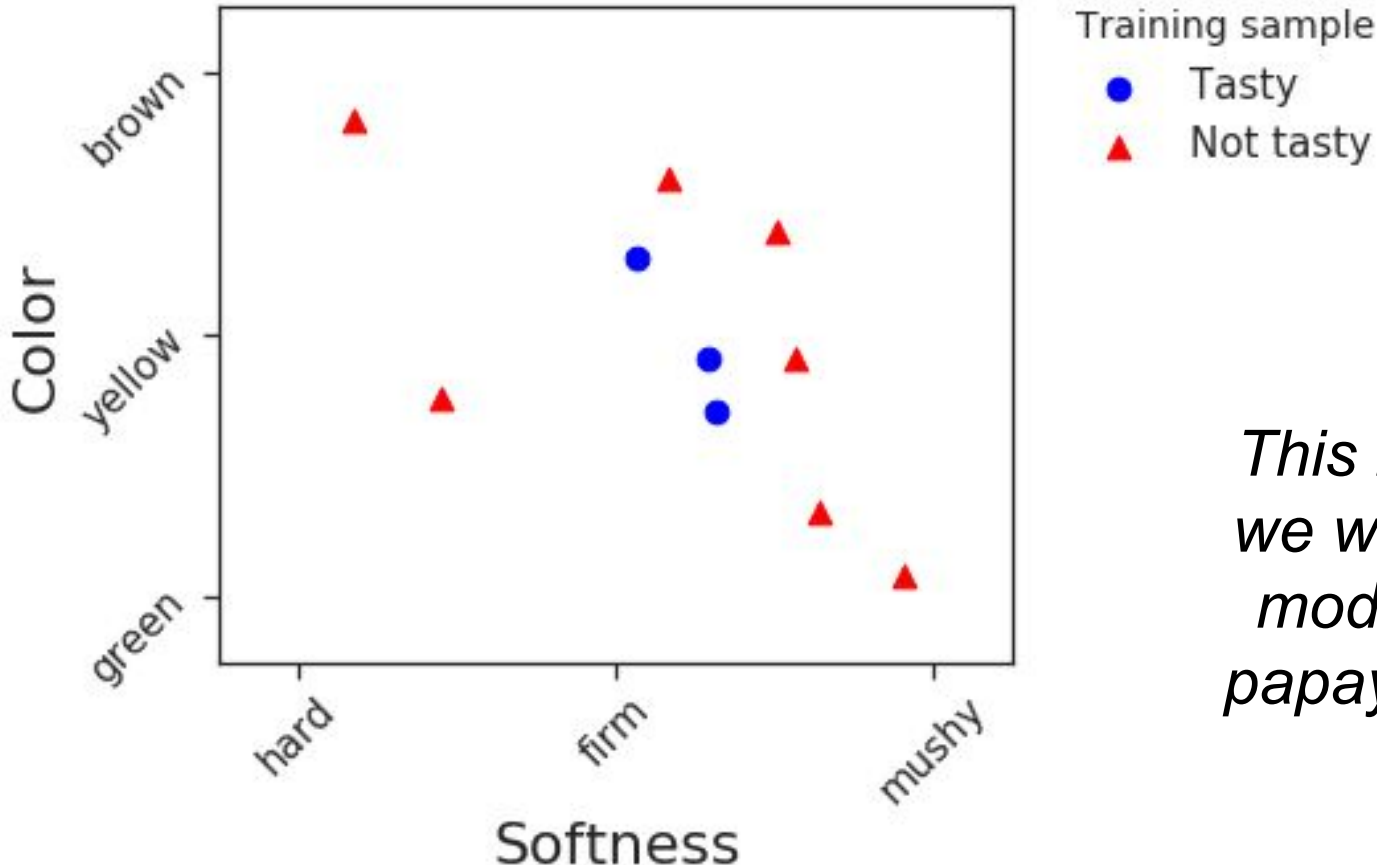
# Take home message

*Priors knowledge is crucial for effective learning*

A controlled example:

# Papaya tasting

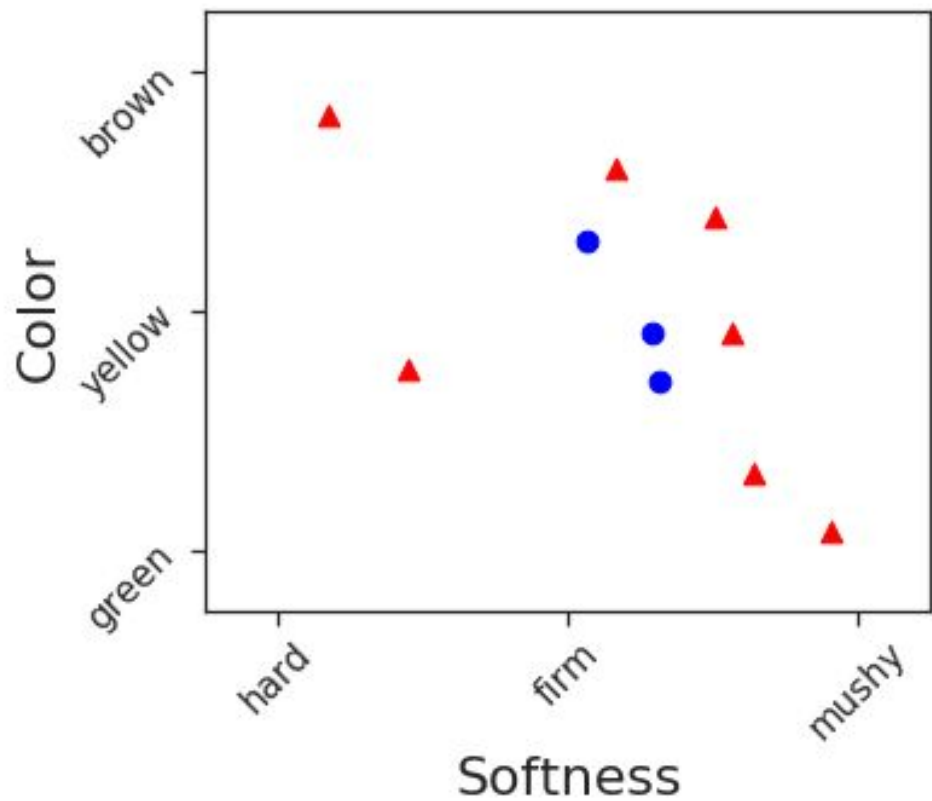Binary classification



*This is all the data we will input to the model about the papayas in the real world!*

YouTube class on the papaya testing example:
https://www.youtube.com/watch?v=b5NlRg8SjZg&list=PLPW2keNyw-usgvmR7FTQ3ZRjfLs5jT4BO&index=2&t=0s

# Papaya tasting



*X:* set of all features,
  *x = [softness, color]*
Y: set of possible labels,
  *y = [tasty, not tasty]*

Training sample
● Tasty
▲ Not tasty

# Papaya tasting



*X:* set of all features,
  *x = [softness, color]*
Y: set of possible labels,
  *y = [tasty, not tasty]*
*D: data generation model,*
  *D ⟹ P(X)*

Training sample
● Tasty
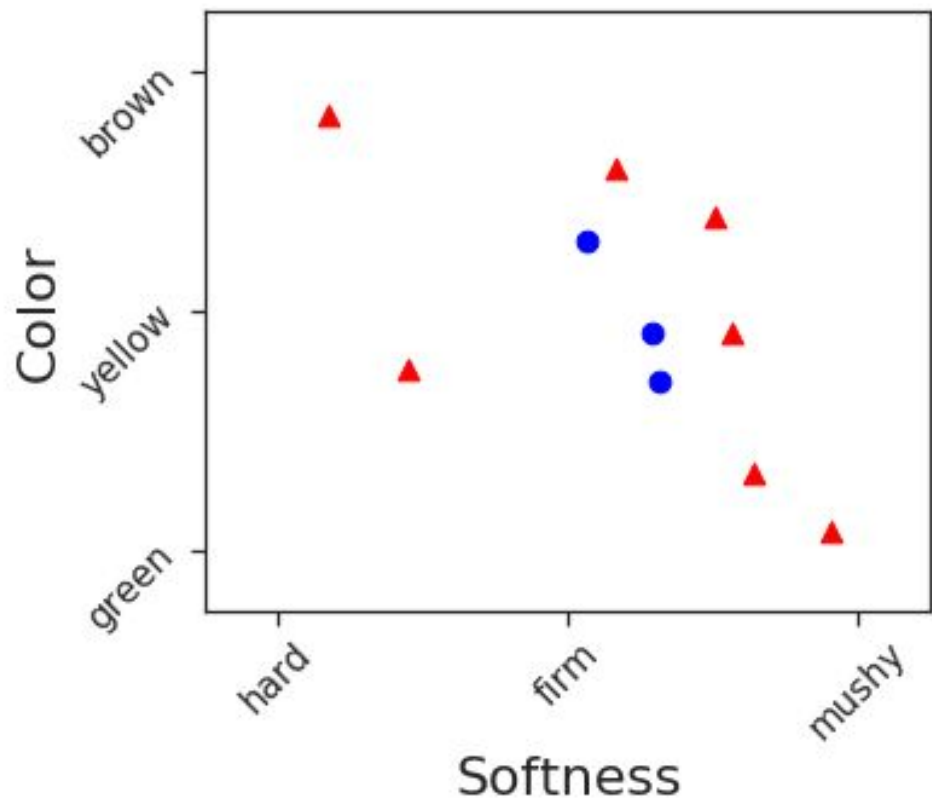▲ Not tasty

A controlled example:

# Papaya tasting



Training sample
- ● Tasty
- ▲ Not tasty

*X:* set of all features,
 *x = [softness, color]*
Y: set of possible labels,
 *y = [tasty, not tasty]*
*D: data generation model,*
 $D \Rightarrow P(X)$
*True Labelling function: y = f(x)*

# Papaya tasting



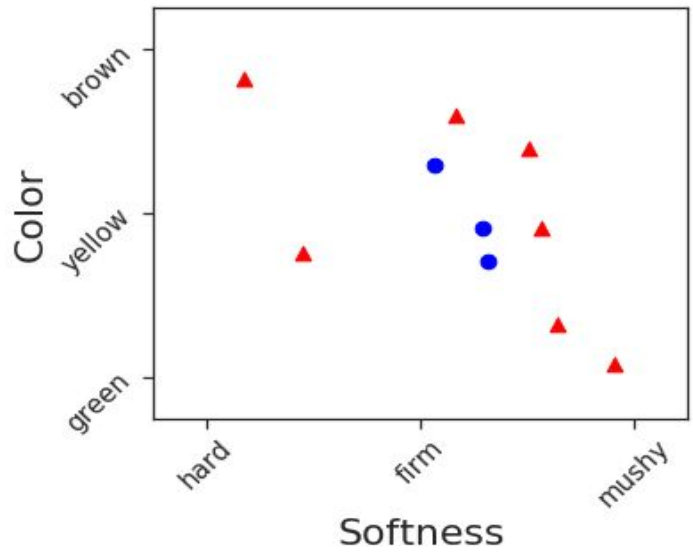*X:* set of all features,
*x = [softness, color]*

Y: set of possible labels,
*y = [tasty, not tasty]*

*D: data generation model,*
$D \Rightarrow P(X)$

*True Labelling function: y = f(x)*

S: training sample: $[x_i, y_i]$, *i ∈ training*

*m: number of objects for training*

Color

brown
yellow
green

hard    firm    mushy

Softness

Training sample
● Tasty
▲ Not tasty

19

# Papaya tasting



$X:$ set of all features,
  $x = [softness, color]$
Y: set of possible labels,
  $y = [tasty, not tasty]$
$D:$ *data generation model,*
  $D \Rightarrow P(X)$
*True Labelling function:* $y = f(x)$

S: training sample: $[x_i, y_i],$ $i \in$ *training*
$m:$ *number of objects for training*
$h_S:$ learner: $y_{est;i} = h_S(x_i)$

Training sample
● Tasty
▲ Not tasty

# Papaya tasting



$X:$ set of all features,
  $x = [softness, color]$
Y: set of possible labels,
  $y = [tasty, not tasty]$
$D:$ *data generation model,*
  $D \Rightarrow P(X)$
*True Labelling function:* $y = f(x)$

S: training sample: $[x_i, y_i],$ *i $\in$ training*
*m: number of objects for training*
$h_S$  learner:  $y_{est;i} = h_S(x_i)$
$L$  metric: $L (y_{true;i} - y_{est;i}),$ *i $\in$ training*

Training sample
● Tasty
▲ Not tasty

Color

brown
yellow
green

Softness
hard
firm
mushy

# Papaya tasting



**Training sample**
- ● Tasty
- ▲ Not tasty

*X:* set of all features,
   *x = [softness, color]*

Y: set of possible labels,
   *y = [tasty, not tasty]*

*D: data generation model,*
   $D \Rightarrow P(X)$

*True Labelling function: y = f(x)*

S: training sample: $[x_i, y_i]$, *i ∈ training*

*m: number of objects for training*

$h_S$    learner:   $y_{est;i} = h_S(x_i)$

L    metric: $L(y_{true;i} - y_{est;i})$, *i ∈ training*

<span style="color:red">Empirical Risk Minimization (ERM)</span>

$L \rightarrow$ *fraction of incorrect predictions*

# Papaya tasting



*X:* set of all features,
   *x = [softness, color]*

Y: set of possible labels,
   *y = [tasty, not tasty]*

*D: data generation model,*
   $D \Rightarrow P(X)$

*True Labelling function: y = f(x)*

S: training sample: *[$x_i$, $y_i$]*, *i ∈ training*

*m: number of objects for training*

$h_S$    learner:   *$y_{est;i}$ = $h_S$($x_i$)*

*L*    metric: *L ($y_{true;i}$ - $y_{est;i}$)*, *i ∈ training*

**Training sample**
- ● Tasty
- ▲ Not tasty

**Empirical Risk Minimization (ERM)**

$$L_{\mathcal{D}}(h_S) = \frac{|\{x \in \mathcal{D} : h_S(x) \neq f(x)\}|}{m}$$
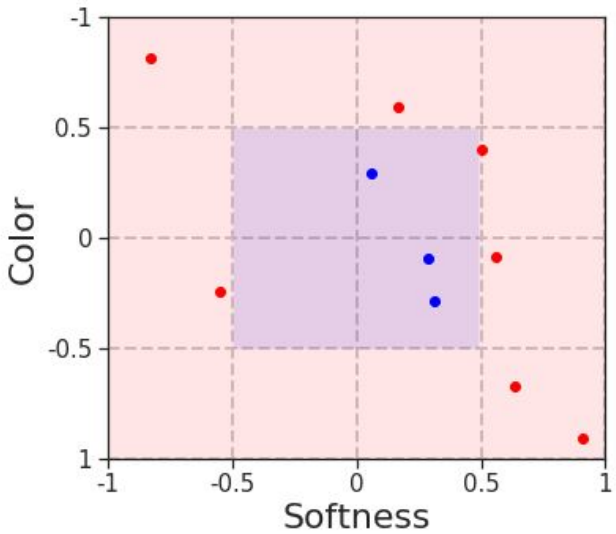
# Papaya tasting

**Proposed learner:**

$$h_S(x) = \begin{cases} y_i & \text{if } x = x_i \mid \{x_i \in S\} \\ 0 & \text{otherwise} \end{cases}$$



Training sample
- Tasty (blue dot)
- Not tasty (red triangle)

*X:* set of all features,
  *x = [softness, color]*
Y: set of possible labels,
  *y = [tasty, not tasty]*
*D: data generation model,*
  $D \Rightarrow P(X)$
*True Labelling function: y = f(x)*

S: training sample: $[x_i, y_i]$, *i ∈ training*
*m: number of objects for training*
$h_S$     learner:   $y_{est;i} = h_S(x_i)$
L     metric: $L (y_{true;i} - y_{est;i})$, *i ∈*
*training*

$$L_{\mathcal{D}}(h_S) = \frac{|\{x \in \mathcal{D} : h_S(x) \neq f(x)\}|}{m}$$
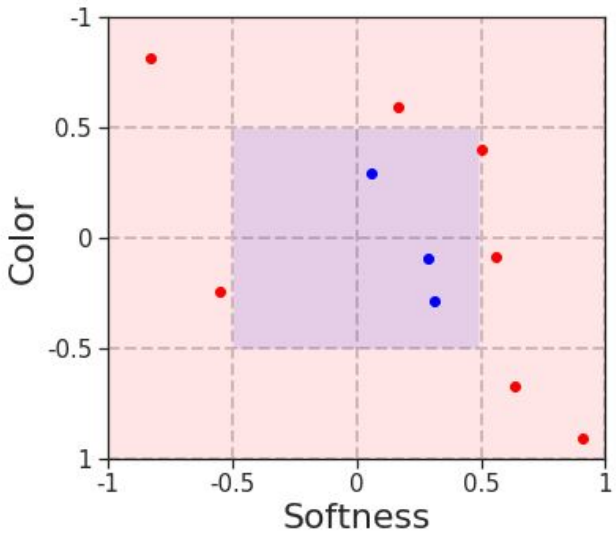
A controlled example:

# Papaya tasting

**Proposed learner:**

$$h_S(x) = \begin{cases} y_i & \text{if } x = x_i \mid \{x_i \in S\} \\ 0 & \text{otherwise} \end{cases}$$

## *Toy model ...*



*X:* set of all features,

   *x = [softness, color]*

Y: set of possible labels,

   *y = [tasty, not tasty]*

*D: data generation model,*

   $D \Rightarrow P(X)$

*True Labelling function: y = f(x)*

S: training sample: *[x_i, y_i],* $i \in$ *training*

*m: number of objects for training*

$h_S$      learner:   *$y_{est;i} = h_S(x_i)$*

$L$      metric: *L ($y_{true;i}$ - $y_{est;i}$),* $i \in$

*training*

$$L_{\mathcal{D}}(h_S) = \frac{|\{x \in \mathcal{D} : h_S(x) \neq f(x)\}|}{m}$$

# Question:

**[tasty, not tasty] = [1, 0]**

## *Proposed learner:*

$$h_S(x) = \begin{cases} y_i & \text{if } x = x_i \mid \{x_i \in S\} \\ 0 & \text{otherwise} \end{cases}$$



Training sample
- Tasty
- Not tasty

What is the expected loss when this model is applied to an arbitrary test sample?

Join at menti.com with code: 7907 6385

*loss* = *fraction of incorrect predictions*

$$L_{\mathcal{D}}(h_S) = \frac{|\{x \in \mathcal{D} : h_S(x) \neq f(x)\}|}{m}$$

26

# Papaya tasting

*Proposed learner:*

$$h_S(x) = \begin{cases} y_i & \text{if } x = x_i \mid \{x_i \in S\} \\ 0 & \text{otherwise} \end{cases}$$

*Answer:*



**Answer:**

$$\begin{aligned} L_S(h_S) &= 0.0 \\ L_D(h_S) &= 0.25 \end{aligned}$$

*X:* set of all features,
  *x = [softness, color]*

Y: set of possible labels,
  **y = [tasty, not tasty] = [1, 0]**

*D: data generation model,*
  *D $\Rightarrow$ P(X)*

*True Labelling function: y = f(x)*

S: training sample: *[$x_i$, $y_i$],* *i $\in$ training*

*m: number of objects for training*

$h_S$   learner: *$y_{est;i} = h_S x_i$*

L   metric: *L ($y_{true,i}$ - $y_{est;i}$),* *i $\in$ training*

$$L_\mathcal{D}(h_S) = \frac{|\{x \in \mathcal{D} : h_S(x) \neq f(x)\}|}{m}$$

# Papaya tasting

*Proposed learner:*

$$h_S(x) = \begin{cases} y_i & \text{if } x = x_i \mid \{x_i \in S\} \\ 0 & \text{otherwise} \end{cases}$$

*Answer:*



$$L_S(h_S) = 0.0$$
$$L_D(h_S) = 0.25$$

*X:* set of all features,
  *x = [softness, color]*
Y: set of possible labels,
  *y = [tasty, not tasty]*
*D: data generation model,*
  $D \Rightarrow P(X)$
*True Labelling function: y = f(x)*

S: training sample: *[x$_i$, y$_i$],* $i \in$ *training*
*m: number of objects for training*
$h_S$   learner:  *y$_{est;i}$ = h$_S$x$_i$*
$L$   metric: *L (y$_{true,i}$ - y$_{est;i}$),* $i \in$
*training*

$$L_{\mathcal{D}}(h_S) = \frac{|\{x \in \mathcal{D} : h_S(x) \neq f(x)\}|}{m}$$

Overfitting! 🥲

28

# Question:

- How can we avoid overfitting?

# Question:

- How can we avoid overfitting?

*by adding prior knowledge ...*

# Choosing the learner

*X:* set of all features,
  *x = [softness, color]*

Y: set of possible labels,
  *y = [tasty, not tasty]*

*D: data generation model,*
  $D \Rightarrow P(X)$

***True Labelling function:*** *y = f(x)*

S: training sample: *[x$_i$, y$_i$],* $i \in training$

*m: number of objects for training*

$h_S$    learner:  $y_{est;i} = h_S(x_i)$

$$h_S(x) = \begin{cases} y_i & \text{if } x = x_i \mid \{x_i \in S\} \\ 0 & \text{otherwise} \end{cases}$$

L: *loss*: *L (y$_{true;i}$ - y$_{est;i}$),* $i \in training$

$$L_{\mathcal{D}}(h_S) = \frac{|\{x \in \mathcal{D} : h_S(x) \neq f(x)\}|}{m}$$

# Choosing the learner

*X:* set of all features,
  *x = [softness, color]*

Y: set of possible labels,
  *y = [tasty, not tasty]*

*D: data generation model,*
  $D \Rightarrow P(X)$

*True Labelling function: y = f(x)*

S: training sample: $[x_i, y_i]$, $i \in training$

*m: number of objects for training*

$h_S$     learner: $y_{est;i} = h_S(x_i)$

$$h_S(x) = \begin{cases} y_i & \text{if } x = x_i \mid \{x_i \in S\} \\ 0 & \text{otherwise} \end{cases}$$

L: *loss*: L $(y_{true;i} - y_{est;i})$, $i \in training$

$$L_{\mathcal{D}}(h_S) = \frac{|\{x \in \mathcal{D} : h_S(x) \neq f(x)\}|}{m}$$

**Hypothesis class (** $\mathcal{H}$ **):**

$$h : \mathcal{X} \longrightarrow \mathcal{Y}; \qquad h \in \mathcal{H}$$

$$\text{ERM}_{\mathcal{H}}(S) \in \underset{h \in \mathcal{H}}{\text{argmin}}\, L_S(h),$$

# Choosing the learner

*X:* set of all features,
  *x = [softness, color]*

Y: set of possible labels,
  *y = [tasty, not tasty]*

*D: data generation model,*
  $D \Rightarrow P(X)$

*True Labelling function: y = f(x)*

S: training sample: *[x<sub>i</sub>, y<sub>i</sub>],* $i \in training$

$h_S$    learner: $y_{est;i} = h_S(x_i)$

$$h_S(x) = \begin{cases} y_i & \text{if } x = x_i \mid \{x_i \in S\} \\ 0 & \text{otherwise} \end{cases}$$

*L:* loss: *L (y<sub>true;i</sub> - y<sub>est;i</sub>),* $i \in training$

$$L_{\mathcal{D}}(h_S) = \frac{|\{x \in \mathcal{D} : h_S(x) \neq f(x)\}|}{m}$$

**Hypothesis class ( $\mathcal{H}$ ):**

$$h : \mathcal{X} \longrightarrow \mathcal{Y}; \qquad h \in \mathcal{H}$$

$$\mathrm{ERM}_{\mathcal{H}}(S) \in \underset{h \in \mathcal{H}}{\arg\min} \, L_S(h),$$

- $\mathcal{H}$ is finite, $N_{\mathcal{H}}$ = number of hypothesis
- The true labelling function is part of $\mathcal{H}$:

$$f \in \mathcal{H}$$

# Choosing the learner

$\chi$: set of all features,
  $x = [softness, color]$

Y: set of possible labels,
  $y = [tasty, not tasty]$

D: *data generation model,*
  $D \Rightarrow P(\chi)$

*True Labelling function: $y = f(x)$*

S: training sample: $[x_i, y_i]$,  $i \in training$

$h_S$     learner:  $y_{est;i} = h_S(x_i)$

$$h_S(x) = \begin{cases} y_i & \text{if } x = x_i \mid \{x_i \in S\} \\ 0 & \text{otherwise} \end{cases}$$

*L: loss*: $L\,(y_{true;i} - y_{est;i})$,  $i \in training$

$$L_{\mathcal{D}}(h_S) = \frac{|\{x \in \mathcal{D} : h_S(x) \neq f(x)\}|}{m}$$

**Hypothesis class ( $\mathcal{H}$ ):**

$$h : \mathcal{X} \longrightarrow \mathcal{Y}; \qquad h \in \mathcal{H}$$

$$\text{ERM}_{\mathcal{H}}(S) \in \operatorname*{argmin}_{h \in \mathcal{H}} L_S(h),$$

- $\mathcal{H}$ is finite, $N_{\mathcal{H}}$ = number of hypothesis
- The true labelling function is part of $\mathcal{H}$:

$$f \in \mathcal{H}$$

- $S$ is identically independently distributed (*i.i.d.*) from $D$

34

# Choosing the learner

$\chi$: set of all features,
$\quad$ x = [softness, color]

Y: set of possible labels,
$\quad$ y = [tasty, not tasty]

D: data generation model,
$\quad$ $D \Rightarrow P(\chi)$

True Labelling function: y = f(x)

S: training sample: [$x_i$, $y_i$], $\;$ i ∈ training

$h_S$ $\quad$ learner: $\;$ $y_{est;i}$ = $h_S(x_i)$

$$h_S(x) = \begin{cases} y_i & \text{if } x = x_i \mid \{x_i \in S\} \\ 0 & \text{otherwise} \end{cases}$$

L: loss: L ($y_{true;i}$ - $y_{est;i}$), $\;$ i ∈ training

$$L_\mathcal{D}(h_S) = \frac{|\{x \in \mathcal{D} : h_S(x) \neq f(x)\}|}{m}$$

**Hypothesis class ( $\mathcal{H}$ ):**

$$h : \mathcal{X} \longrightarrow \mathcal{Y}; \qquad h \in \mathcal{H}$$

$$\text{ERM}_\mathcal{H}(S) \in \underset{h \in \mathcal{H}}{\arg\min}\, L_S(h),$$

- $\mathcal{H}$ is finite, $N_\mathcal{H}$ = number of hypothesis
- The true labelling function is part of $\mathcal{H}$:

$$f \in \mathcal{H}$$

- S is identically independently distributed (i.i.d.) from D

35

# Machine Learning:

*(a personal favorite)*
## Supervised definition

Hypothesis: x $\longrightarrow$ | Nature | $\longrightarrow$ y

*Breiman, L., Statistical Modeling: The Two Cultures, Stat. Sci, Volume 16 (2001)*

*Breiman, L., Statistical Modeling: The Two Cultures, Stat. Sci, Volume 16 (2001)*

**Hypothesis:**

x → Nature → y

**Physical modeling:**

x → [ Linear regression / Logistic regression / … / **Statistical model** ] → y

**Algorithmic modeling:**

x → Nature → y

[ Nearest Neighbor / Decision trees / … / **Machine Learning model** ]

*Breiman, L., Statistical Modeling: The Two Cultures, Stat. Sci, Volume 16 (2001)*

**Hypothesis:** x → Nature → y

**Physical modeling:** x → [ Linear regression / Logistic regression / … / Statistical model ] → y

**Algorithmic modeling:** x → Nature → y; Why should this work? [ Nearest Neighbor / Decision trees / … / Machine Learning model ]

*Breiman, L., Statistical Modeling: The Two Cultures, Stat. Sci, Volume 16 (2001)*

# Representativeness

Probability distribution, *P*



$$(\mu_P, \sigma_P)$$

Sample, S1



$$(\mu_{S_1}, \sigma_{S_1})$$

# Representativeness

Probability distribution, *P*



$$(\mu_P, \sigma_P)$$

Sample, S1



$$(\mu_{S_1}, \sigma_{S_1})$$

# Representativeness

Probability distribution, $P$



$$(\mu_P, \sigma_P)$$

$S_1$ is *representative* of $P$

Sample, S1



$$(\mu_{S_1}, \sigma_{S_1})$$

43

# Representativeness

Probability distribution, $P$

$(\mu_P, \sigma_P)$



f(x)

$\sigma$

$\mu$

x

Sample, $S_1$

Sample, $S_2$

$(\mu_{S_1}, \sigma_{S_1})$

$(\mu_{S_2}, \sigma_{S_2})$

# Representativeness

Probability distribution, *P*



$(\mu_P, \sigma_P)$

Sample, $S_1$

Sample, $S_2$

$(\mu_{S_1}, \sigma_{S_1})$ $(\mu_{S_2}, \sigma_{S_2})$

# Representativeness

Probability distribution, *P*



$$(\mu_P, \sigma_P)$$

This is why it works!

Training

Test

$$(\mu_{S_1}, \sigma_{S_1})$$       $$(\mu_{S_2}, \sigma_{S_2})$$

# Representativeness

- A sample S1 is said to be representative of a probability distribution P if one can draw accurate conclusions about P from S1

- If two samples S1 and S2 are representative of P, S1 and S2 are representative in relation to each other

# Representativeness

- A sample S1 is said to be representative of a probability distribution P if one can draw accurate conclusions about P from S1

- If two samples S1 and S2 are representative of P, S1 and S2 are representative in relation to each other

## Question:

*If a sample $S_1$ identically independently distributed (i.i.d.) from a distribution P, is this enough to guarantee that $S_1$ is representative of P?*

# Model assumptions

$\chi$:  set of all features,
  $x = [softness, color]$

Y:  set of possible labels,
  $y = [tasty, not tasty]$

D: data generation model,
  $D \Rightarrow P(\chi)$

True Labelling function: $y = f(x)$

S:  training sample: $[x_i, y_i]$,  $i \in training$

$h_S$    learner:  $y_{est;i} = h_S(x_i)$

$$h_S(x) = \begin{cases} y_i & \text{if } \exists i \in [m] \text{ s.t. } x_i = x \\ 0 & \text{otherwise.} \end{cases}$$

L:  loss: $L(y_{true;i} - y_{est;i})$,  $i \in training$

$$L_{\mathcal{D}}(h_S) = \frac{|\{x \in \mathcal{D} : h_S(x) \neq f(x)\}|}{m}$$

**Hypothesis class ( $\mathcal{H}$ ):**

$$h : \mathcal{X} \longrightarrow \mathcal{Y}; \qquad h \in \mathcal{H}$$

$$\text{ERM}_{\mathcal{H}}(S) \in \underset{h \in \mathcal{H}}{\arg\min}\, L_S(h),$$

- $\mathcal{H}$ is finite, $N_{\mathcal{H}}$ = number of hypothesis

- The true labelling function is part of $\mathcal{H}$:

  $$f \in \mathcal{H}$$

- $S$ is identically independently distributed (*i.i.d.*) from $D$

*Not enough!*

# Bad samples and hypothesis

# Bad samples and hypothesis

$\delta \rightarrow$ probability of non-representative (bad) samples

# Bad hypothesis and samples

$\delta \rightarrow$ probability of non-representative (bad) samples

$1 - \delta \rightarrow$ confidence parameter

# Bad hypothesis and samples

$\delta \rightarrow$ probability of non-representative (bad) samples

$1 - \delta \rightarrow$ confidence parameter

$\varepsilon \rightarrow$ contamination. A failure will occur when $\quad L_D(h_S) \quad \geq \quad \epsilon$

Good hypothesis:
$$\mathcal{H}_G \ := \ [h \in \mathcal{H} : L_S(h_S) = 0 \quad \& \quad L_D(h_S) < \epsilon]$$

Bad hypothesis:
$$\mathcal{H}_B \ := \ [h \in \mathcal{H} : L_S(h_S) = 0 \quad \& \quad L_D(h_S) \geq \epsilon]$$

# Bad hypothesis and samples

$\delta \rightarrow$ probability of non-representative (bad) samples

$1 - \delta \rightarrow$ confidence parameter

$\varepsilon \rightarrow$ contamination. A failure will occur when $\quad L_D(h_S) \quad \geq \quad \epsilon$

| | | | | |
|---|---|---|---|---|
| Good hypothesis: | $\mathcal{H}_G$ | $:=$ | $[h \in \mathcal{H} : L_S(h_S) = 0$ | $\&\quad L_D(h_S) < \epsilon]$ |
| Bad hypothesis: | $\mathcal{H}_B$ | $:=$ | $[h \in \mathcal{H} : L_S(h_S) = 0$ | $\&\quad L_D(h_S) \geq \epsilon]$ |

Realizability assumption, $f \in \mathcal{H}$

# Constructing misleading samples

*The world*



Bad samples

*For 1 element in the training sample*

$$x_i \quad | \quad h(x_i) \quad = \quad y_i$$

# Constructing misleading samples

*The world*



Bad samples

*For 1 element in the training sample*

$$x_i \quad | \quad h(x_i) \quad = \quad y_i$$

$$P(x_i \in \mathcal{D} : h(x_i) \quad = \quad y_i) = 1 - L_{\mathcal{D},f}(h)$$

# Constructing misleading samples

**The world**



Bad samples

*For 1 element in the training sample*

$$x_i \quad | \quad h(x_i) \quad = \quad y_i$$

$$P(x_i \in \mathcal{D} : h(x_i) \quad = \quad y_i) = 1 - L_{\mathcal{D},f}(h)$$

$$P(x_i \in \mathcal{D} : h(x_i) \quad = \quad y_i) \le 1 - \epsilon$$

Things can still go wrong …

# Constructing misleading samples

*The world*



Bad samples

*For 1 element in the training sample*

$$x_i \quad | \quad h(x_i) \quad = \quad y_i$$

$$P(x_i \in \mathcal{D} : h(x_i) \quad = \quad y_i) = 1 - L_{\mathcal{D},f}(h)$$

$$P(x_i \in \mathcal{D} : h(x_i) \quad = \quad y_i) \leq 1 - \epsilon$$

*For **m** elements in the training sample*

Since all elements in training are i.i.d.,

$$P(S_m : L_S(h) = 0) \quad \leq \quad \prod_{i=1}^{m}(1 - \epsilon) = (1 - \epsilon)^m$$

# Considering bad hypothesis

*For **1** hypothesis*

$$P(S_m : L_S(h) = 0) \quad \leq \quad (1 - \epsilon)^m$$

# Considering bad hypothesis



*For **1** hypothesis*

$$P(S_m : L_S(h) = 0) \quad \leq \quad (1 - \epsilon)^m$$

The sum rule

$$P(A \cup B) \quad \leq \quad P(A) + P(B)$$

# Considering bad hypothesis

*For **1** hypothesis*

$$P(S_m : L_S(h) = 0) \quad \leq \quad (1 - \epsilon)^m$$

The sum rule

$$P(A \cup B) \quad \leq \quad P(A) + P(B)$$

*For all bad hypothesis*

$$\delta = P(L_S(h) = 0, \forall h \in \mathcal{H}_B) \quad \leq \quad \sum_{h \in \mathcal{H}_B} (1 - \epsilon)^m$$

# Considering bad hypothesis



*For **1** hypothesis*

$$P(S_m : L_S(h) = 0) \quad \leq \quad (1 - \epsilon)^m$$

The sum rule

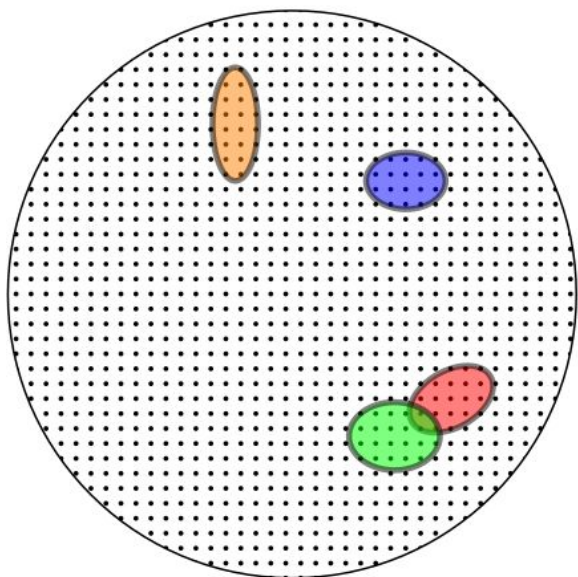$$P(A \cup B) \quad \leq \quad P(A) + P(B)$$

*For all bad hypothesis*

$$\delta = P(L_S(h) = 0, \forall h \in \mathcal{H}_B) \quad \leq \quad \sum_{h \in \mathcal{H}_B} (1 - \epsilon)^m$$

using…

$$(1 - x)^y \quad \leq \quad \exp(-xy)$$

# Considering bad hypothesis

*For **1** hypothesis*

$$P(S_m : L_S(h) = 0) \quad \leq \quad (1 - \epsilon)^m$$

The sum rule

$$P(A \cup B) \quad \leq \quad P(A) + P(B)$$

*For all bad hypothesis*

$$\delta = P(L_S(h) = 0, \forall h \in \mathcal{H}_B) \quad \leq \quad \sum_{h \in \mathcal{H}_B} (1 - \epsilon)^m$$

using...

$$(1 - x)^y \quad \leq \quad \exp(-xy)$$

$$\boxed{\delta \leq N_{\mathcal{H}} \exp(-\epsilon m)}$$

# In summary …

$$\delta \leq N_{\mathcal{H}} \exp(-\epsilon m)$$

*L. Valiant. A theory of the learnable. Communications of the ACM, 27, 1984.*

# In summary …

$$\delta \leq N_{\mathcal{H}} \exp(-\epsilon m)$$

If,
$$m_{\mathcal{H}}(\epsilon, \delta) \geq \frac{\ln(N_{\mathcal{H}}/\delta)}{\epsilon} \longrightarrow$$
every $h$ from ERM, $L_{(\mathcal{D}, f)}(h_S) \leq \epsilon.$

*L. Valiant. A theory of the learnable. Communications of the ACM, 27, 1984.*

# PAC learning model

$$\delta \leq N_{\mathcal{H}} \exp(-\epsilon m)$$

If,                                                           every *h* from ERM,

$$m_{\mathcal{H}}(\epsilon, \delta) \geq \frac{\ln(N_{\mathcal{H}}/\delta)}{\epsilon} \longrightarrow L_{(\mathcal{D},f)}(h_S) \leq \epsilon.$$

*L. Valiant. A theory of the learnable. Communications of the ACM, 27, 1984.*

# PAC learning model

$$\delta \leq N_{\mathcal{H}} \exp(-\epsilon m)$$

**P**robably $\qquad \rightarrow$ with confidence 1 - δ over *m* samples

**A**pproximately $\quad \rightarrow$ within a contamination level $\leq$ ε

**C**orrect

If, $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ every *h* from ERM,

$$m_{\mathcal{H}}(\epsilon, \delta) \geq \frac{\ln(N_{\mathcal{H}}/\delta)}{\epsilon} \quad\longrightarrow\quad L_{(\mathcal{D},f)}(h_S) \leq \epsilon.$$

*L. Valiant. A theory of the learnable. Communications of the ACM, 27, 1984.*

# PAC Assumptions

$\chi$: set of all features,
  $x = [softness, color]$

Y: set of possible labels,
  $y = [tasty, not\ tasty]$

D: data generation model,
  $D \Rightarrow P(\chi)$

True Labelling function: $y = f(x)$

S: training sample: $[x_i, y_i]$, $i \in training$

m: number of objects for training

$h_S$    learner: $y_{est;i} = h_S(x_i)$

$$h_S(x) = \begin{cases} y_i & \text{if } \exists i \in [m] \text{ s.t. } x_i = x \\ 0 & \text{otherwise.} \end{cases}$$

L: loss: $L(y_{true;i} - y_{est;i})$, $i \in training$

$$L_{\mathcal{D}}(h_S) = \frac{|\{x \in \mathcal{D} : h_S(x) \neq f(x)\}|}{m}$$

**Hypothesis class ( $\mathcal{H}$ ):**

$$h : \mathcal{X} \longrightarrow \mathcal{Y}; \qquad h \in \mathcal{H}$$

$$\text{ERM}_{\mathcal{H}}(S) \in \underset{h \in \mathcal{H}}{\arg\min}\ L_S(h),$$

- $\mathcal{H}$ is finite, $N_{\mathcal{H}}$= number of hypothesis

- The true labelling function is part of $\mathcal{H}$:

$$f \in \mathcal{H}$$

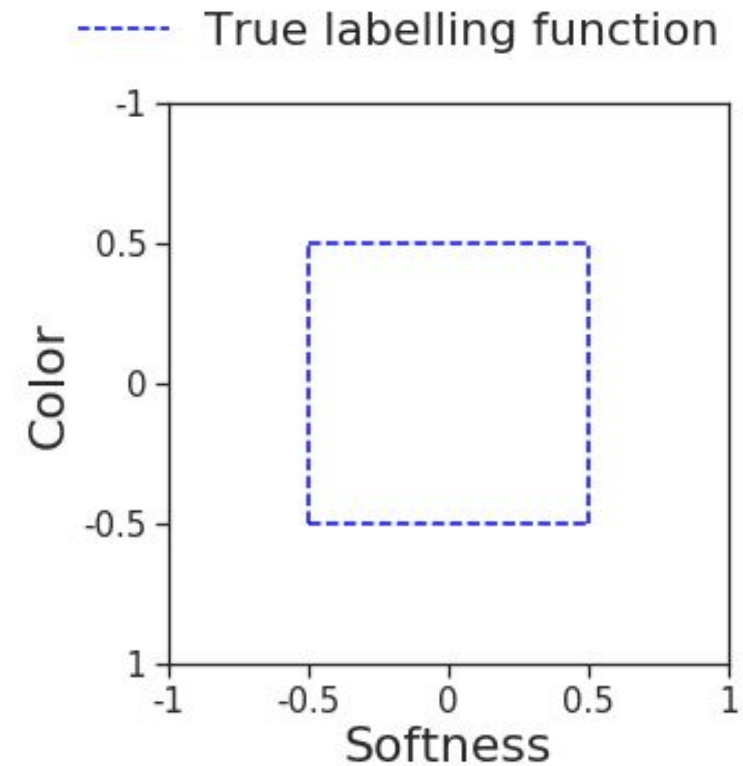- S is identically independently distributed (i.i.d.) from D
- Representativeness

# Papaya tasting



$\chi$:  set of $x \in$ [softness, color]

Y:  set  y = [tasty, not tasty]

D: data generation model:  $D \Rightarrow P(\chi)$

# Papaya tasting

*χ:*  set of  *x* ∈ *[softness, color]*
Y:  set  *y = [tasty, not tasty]*
*D: data generation model:  D ⇒ P(χ)*
**True Labelling function:**
*y = tasty if softness* ∈ *[-0.5, 0.5] and*
              *color* ∈ *[-0.5, 0.5]*



----- True labelling function

# Papaya tasting

$\chi$:  set of $x \in$ [softness, color]

Y: set $y$ = [tasty, not tasty]

*D: data generation model: $D \Rightarrow P(\chi)$*

***True Labelling function:***

*y = tasty if softness $\in$ [-0.5, 0.5] and*
                    *color $\in$ [-0.5, 0.5]*

S:  training sample: $[x_i, y_i]$, $i \in$ training

*m: number of objects for training*



----- True labelling function

# Papaya tasting

$\chi$:  set of $x \in$ [softness, color]

Y:  set  y = [tasty, not tasty]

D: data generation model:  $D \Rightarrow P(\chi)$

**True Labelling function:**

y = tasty if softness $\in$ [-0.5, 0.5] and
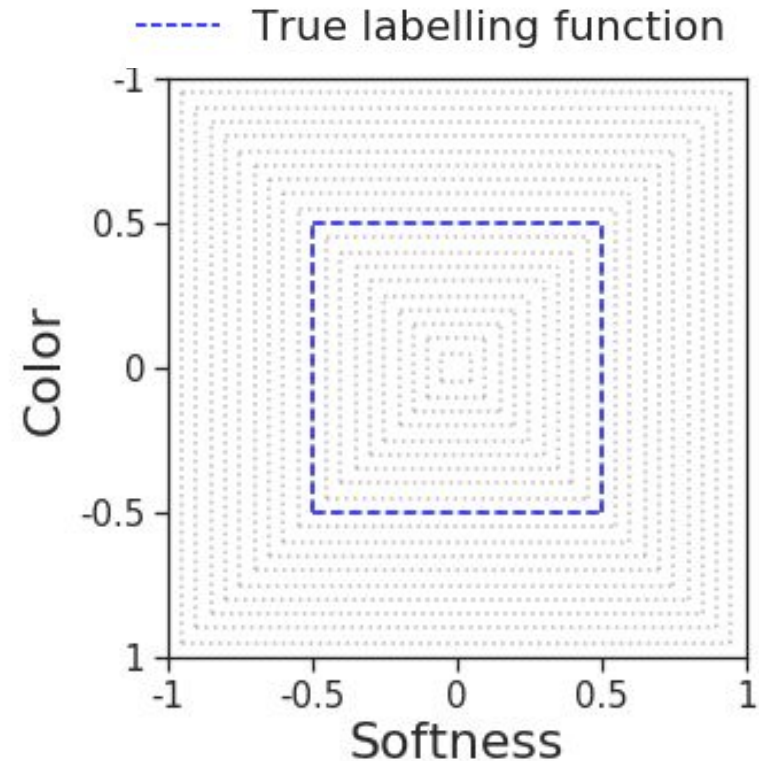                color $\in$ [-0.5, 0.5]

S:  training sample: [$x_i$, $y_i$],  $i \in training$

m: number of objects for training

$\mathcal{H}$: hypothesis class:

   axis aligned squares in steps of 0.05

$N_H$ = 20

# Papaya tasting

$\chi$:  set of $x \in$ [softness, color]

Y:  set  $y$ = [tasty, not tasty]

D: data generation model:  $D \Rightarrow P(\chi)$

**True Labelling function:**

$y$ = tasty if softness $\in$ [-0.5, 0.5] and

color $\in$ [-0.5, 0.5]

S:  training sample: $[x_i, y_i]$, $i \in training$
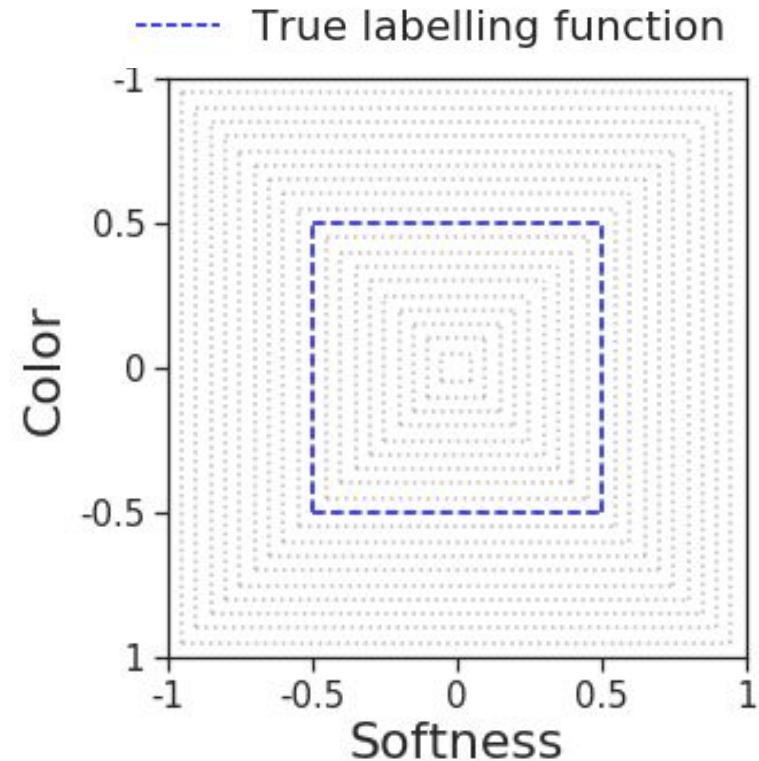
m: number of objects for training

$\mathcal{H}$: hypothesis class:

axis aligned squares in steps of 0.05

$N_H$ = 20

L:  loss: $L (y_{true;i} - y_{est;i})$, $i \in training$

$$L_{\mathcal{D}}(h_S) = \frac{|\{x \in \mathcal{D} : h_S(x) \neq f(x)\}|}{m}$$



True labelling function

# Question:



Data model: uniform distribution [-1,1] in both axis

$1 - \delta = 0.95 \leftarrow$    confidence

$\varepsilon \quad = 0.05 \leftarrow$    contamination

$N_H = 20 \quad \leftarrow$    number of possible squares

**m = ??**

Join at menti.com with code: 7907 6385

*What would you guess is the number of examples necessary for training?*

# Question:

Data model: uniform distribution [-1,1] in both axis

$1 - \delta = 0.95 \leftarrow$ confidence

$\varepsilon = 0.05 \leftarrow$ contamination

$N_H = 20 \leftarrow$ number of possible squares

**m ~ 120**



----- True labelling function

# Question:

Data model: uniform distribution
[-1,1] in both axis

$1 - \delta = 0.95 \leftarrow$ confidence
$\varepsilon \quad = 0.05 \leftarrow$ contamination
$N_H = 20 \quad \leftarrow$ number of possible squares



Color / Softness scatter plot with legend: ---- True labelling function



minimum training vs number of hypothesis plot with legend:

$\delta = 0.01, \varepsilon = 0.01$
$\delta = 0.05, \varepsilon = 0.01$
$\delta = 0.01, \varepsilon = 0.05$
$\delta = 0.05, \varepsilon = 0.05$

# Agnostic PAC learning

$\chi$: set of all features,

   $x = [softness, color]$

Y: set of possible labels,

   $y = [tasty, not\ tasty]$

D: *data generation model,*

   $D \Rightarrow P(\chi, Y)$

*True Labelling function:* $y = f([x,y])$

S: training sample: $[x_i, y_i]$, $i \in training$

*m: number of objects for training*

$h_S$     learner:    $y_{est;i} = h_S(x_i, y_i)$

*L: loss*

$$L_{\mathcal{D}}(h) \stackrel{\text{def}}{=} \mathop{\mathbb{E}}_{(x,y)\sim\mathcal{D}} (h(x) - y)^2$$

$$h : \mathcal{X} \longrightarrow \mathcal{Y}; \qquad h \in \mathcal{H}$$

$$\text{ERM}_{\mathcal{H}}(S) \in \operatorname*{argmin}_{h\in\mathcal{H}} L_S(h),$$

- $m \rightarrow$ number of objects in training

- $\mathcal{H}$ is finite, $N_H$ = number of hypothesis

- The true labelling function may not be part of $\mathcal{H}$:

$$f \notin \mathcal{H}$$

$$L_{\mathcal{D}}(h) \leq \min_{h'\in\mathcal{H}} L_{\mathcal{D}}(h') + \epsilon,$$

# Representativeness
## *in machine learning*

*Shalev-Shwartz, S. and Ben-David, S., Understanding Machine Learning - from theory to algorithms, 2014, Cambridge University Press*

# Representativeness

## *in machine learning*

### *or*

## **Uniform Convergence**

$$\forall h \in \mathcal{H}, \quad |L_S(h) - L_{\mathcal{D}}(h)| \leq \epsilon$$

# Representativeness
*in machine learning*

## *or*

## ***Uniform Convergence***

$$\forall h \in \mathcal{H}, \quad |L_S(h) - L_{\mathcal{D}}(h)| \leq \epsilon$$

<u>It can be shown that,</u> if $\mathcal{H}$ has uniform convergence, ERM$_{\mathcal{H}}$ is a successful agnostic PAC learner of $\mathcal{H}$.

# Can machine learning solve my problem?

# Can machine learning solve my problem?

- If your data satisfy all the necessary conditions;

# Can machine learning solve my problem?

- If your data satisfy all the necessary conditions;

- If you have enough training sample to fulfill your expectations;

# Can machine learning solve my problem?

- If your data satisfy all the necessary conditions;

- If you have enough training sample to fulfill your expectations;

- If the set of your hypothesis class + loss function + training data has uniform convergence (representativeness)

# Can machine learning solve my problem?

- If your data satisfy all the necessary conditions;

- If you have enough training sample to fulfill your expectations;

- If the set of your hypothesis class + loss function + training data has uniform convergence (representativeness)

*Then.. probably (1-δ), approximately (ε) : yes*

# What about practical situations?

*If you are using a classical learner whose class under your training sample and loss function are representative (has uniform convergence), you are probably getting reasonable results… **but not all the time**!*

# What about practical situations?

*If you are using a classical learner whose class under your training sample and loss function are representative (has uniform convergence), you are probably getting reasonable results… **but not all the time***!

So why does it seem to work in everything around us?

Best guess:  *we do not know how to model real data…*

*https://www.youtube.com/watch?v=M2BJC0OyI-w*

# There is plenty room for improvement!

*Progress will only be possible through **interdisciplinary** collaboration!*

# There is plenty room for improvement!

*Progress will only be possible through **interdisciplinary** collaboration!*

Machine learning is a wonderful field of research, which has already shown its potential in many fields! We should definitely take advantage of its results .. however ...
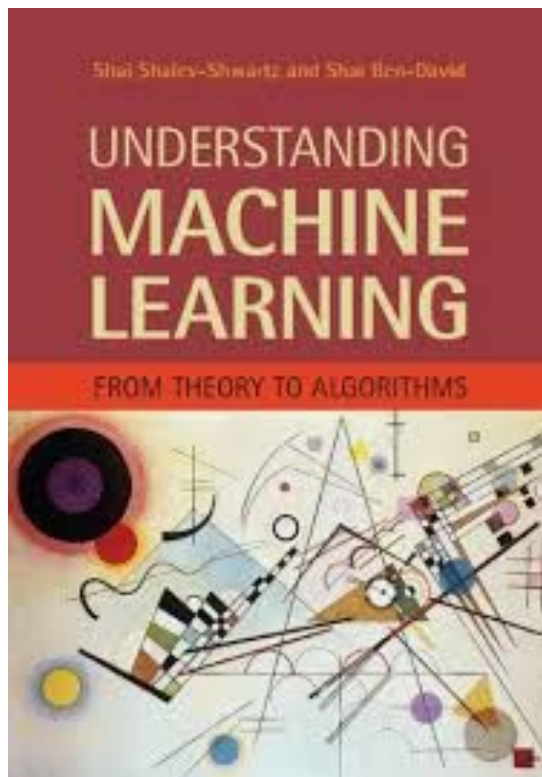
# There is plenty room for improvement!

*Progress will only be possible through **interdisciplinary** collaboration!*

Machine learning is a wonderful field of research, which has already shown its potential in many fields! We should definitely take advantage of its results .. however ...

*References:*

# This talk is a rough summary of chapters 1-4:

*Free download - with agreement from the editor:*

https://www.cse.huji.ac.il/~shais/UnderstandingMachineLearning/index.html

*23 lectures of 1.5 hours each on youtube:*

https://www.youtube.com/playlist?list=PLPW2keNyw-usgvmR7FTQ3ZRjfLs5jT4BO

*Enjoy!*

THANK YOU