

Introduction to Machine Learning Theory II

ML@Cezeaux, 18 February 2020

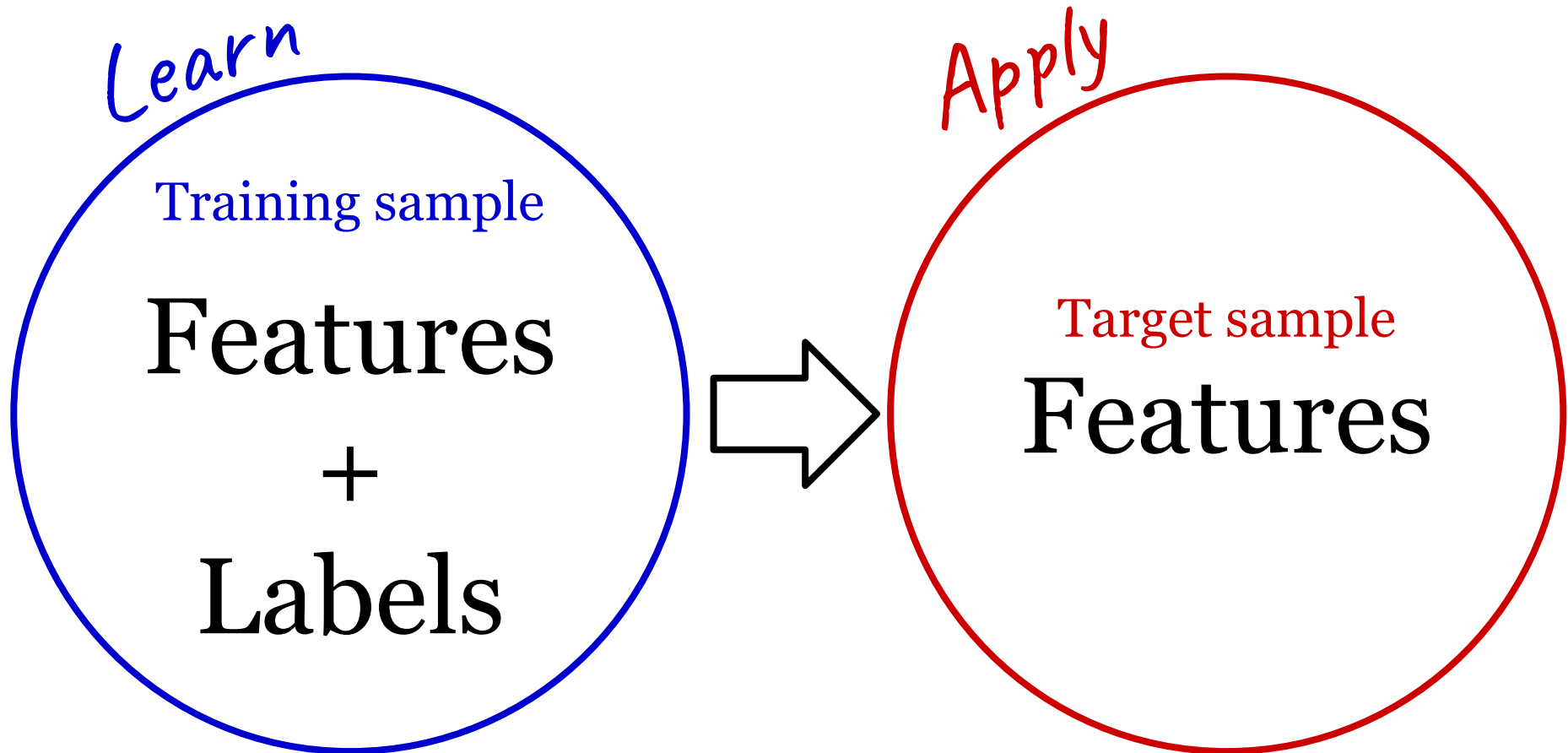
Emille E. O. Ishida

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

Start from the beginning ...

Supervised Learning

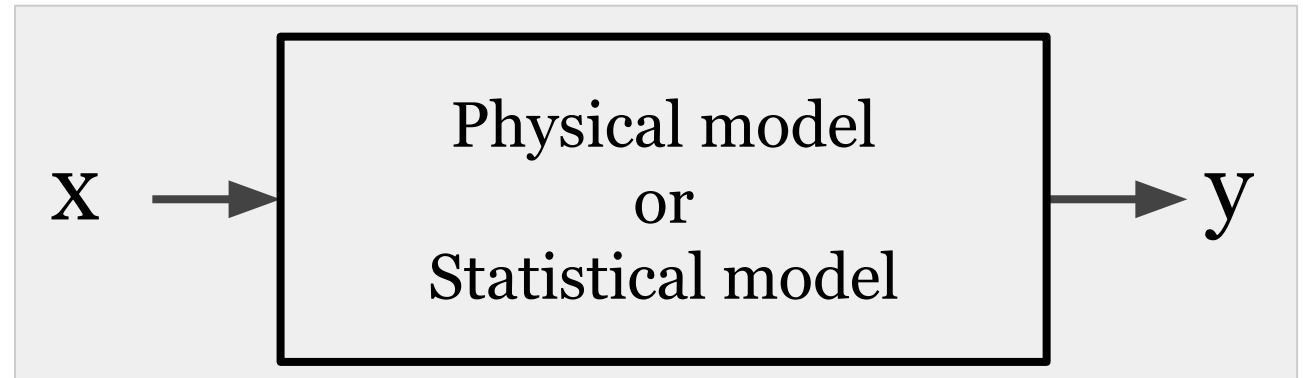
Learn by example



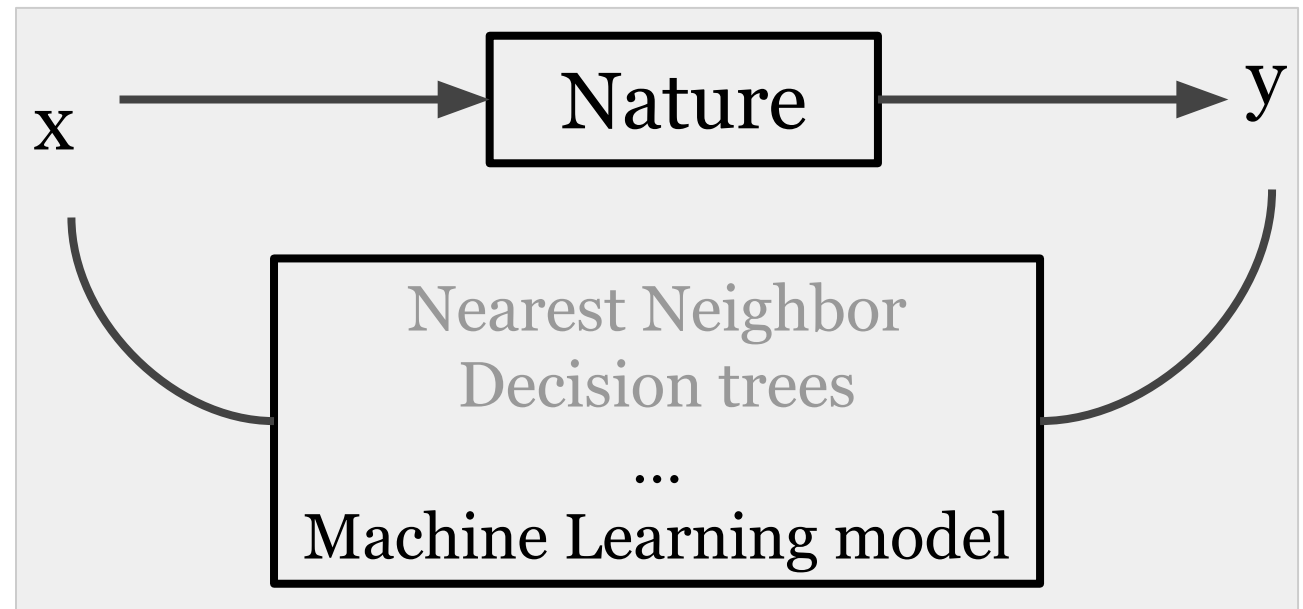
Hypothesis:



Physical modeling:

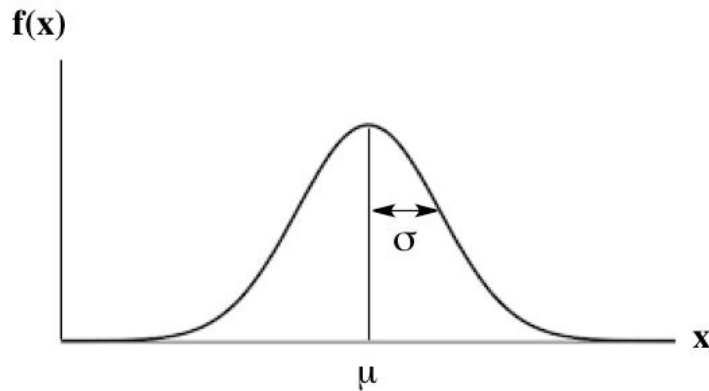


Algorithmic modeling:



Representativeness

Probability distribution, P



$$(\mu_P, \sigma_P)$$

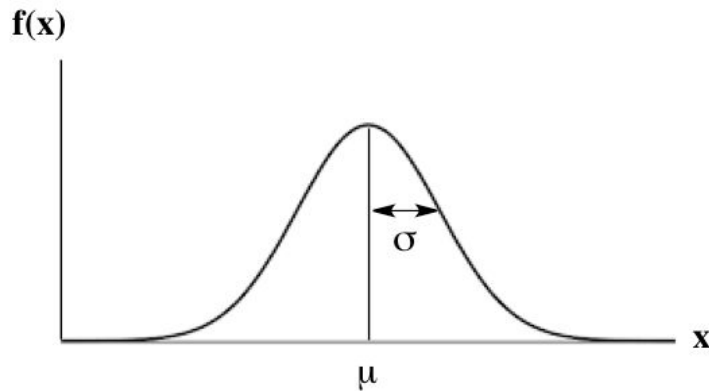
Sample, S_1



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

Representativeness

Probability distribution, P



$$(\mu_P, \sigma_P)$$



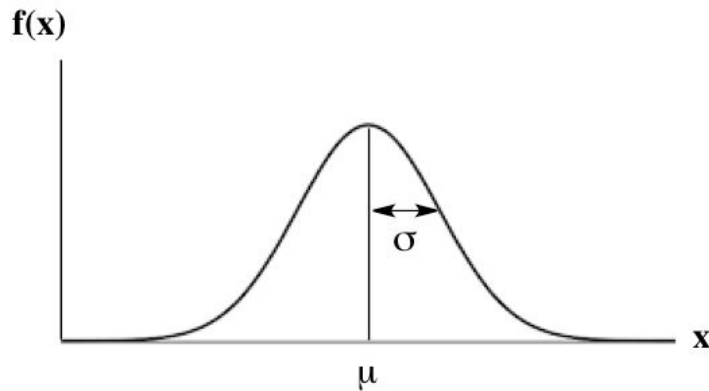
Sample, S_1



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

Representativeness

Probability distribution, P



$$(\mu_P, \sigma_P)$$

Sample, S_1

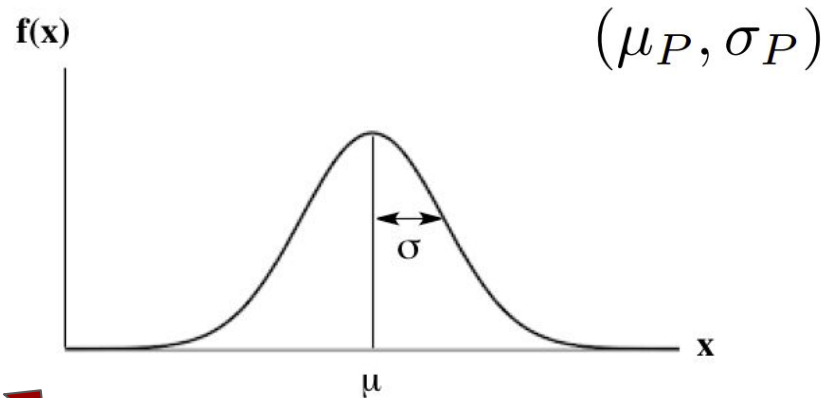


$$(\mu_{S_1}, \sigma_{S_1})$$

S_1 is
representative
of P

Representativeness

Probability distribution, P



Sample, S_1



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

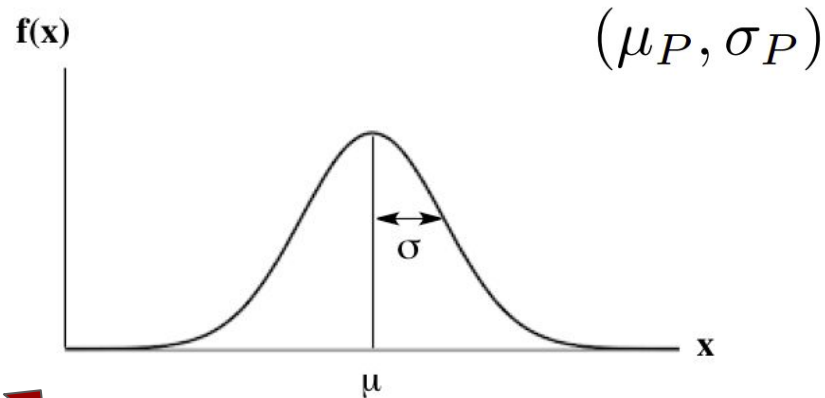
Sample, S_2



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

Representativeness

Probability distribution, P



Sample, S_1



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

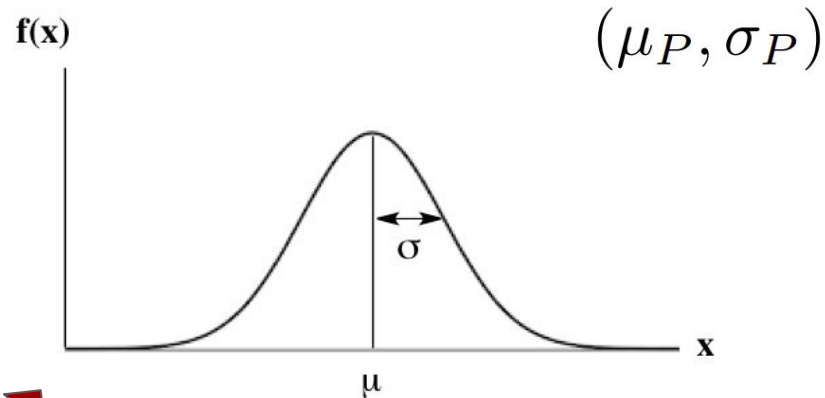
Sample, S_2



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

Representativeness

Probability distribution, P



Training
Features
+
Labels



Target
Features
+
(?)

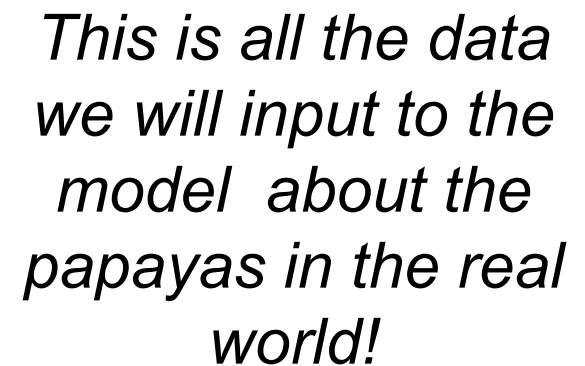


$(\mu_{S_1}, \sigma_{S_1})$

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



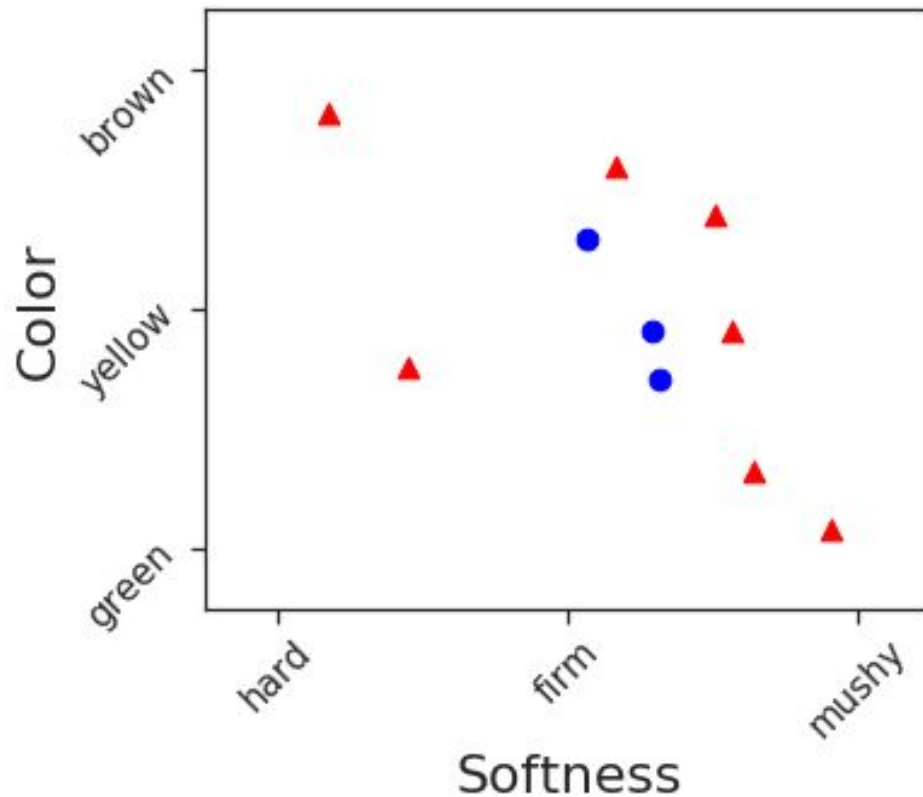
Binary classification



<https://www.youtube.com/watch?v=b5NIRg8SiZg&list=PLPW2keNyw-usqvmR7FTQ3ZRifLs5iT4BO&index=2&t=0s>

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: $y = f(x)$

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

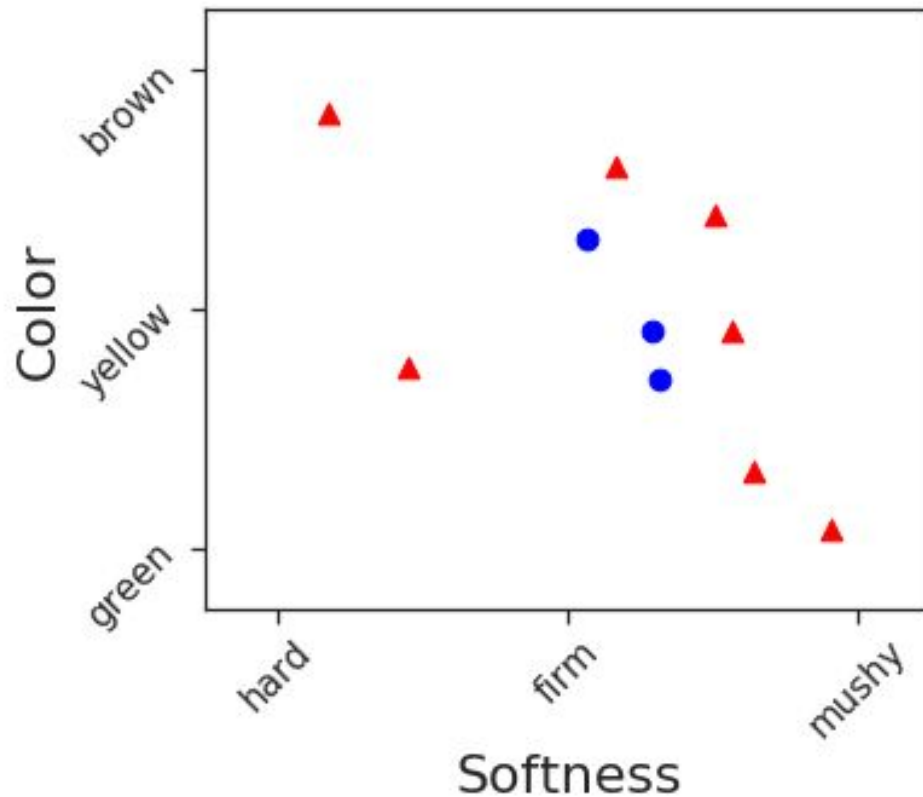
m: number of objects for training

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

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

A controlled example:

Papaya tasting



Training sample

- Tasty
- ▲ Not tasty

Empirical Risk Minimization (ERM)

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: $[x_i, y_i], i \in \text{training}$

m : number of objects for training

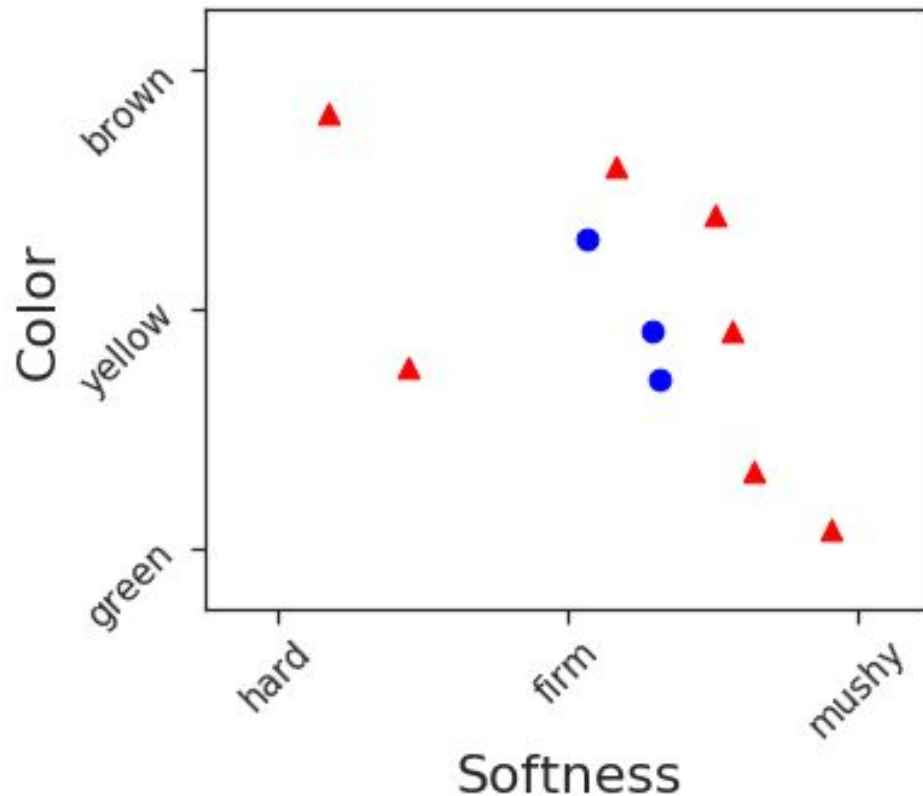
$$h_S \quad \text{learner: } y_{est,i} = h_S(x_i)$$

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

$L \rightarrow$ fraction of incorrect predictions

A controlled example:

Papaya tasting



Training sample

- Tasty
- ▲ Not tasty

Empirical Risk Minimization (ERM)

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: $[x_i, y_i], i \in \text{training}$

m : number of objects for training

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

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

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

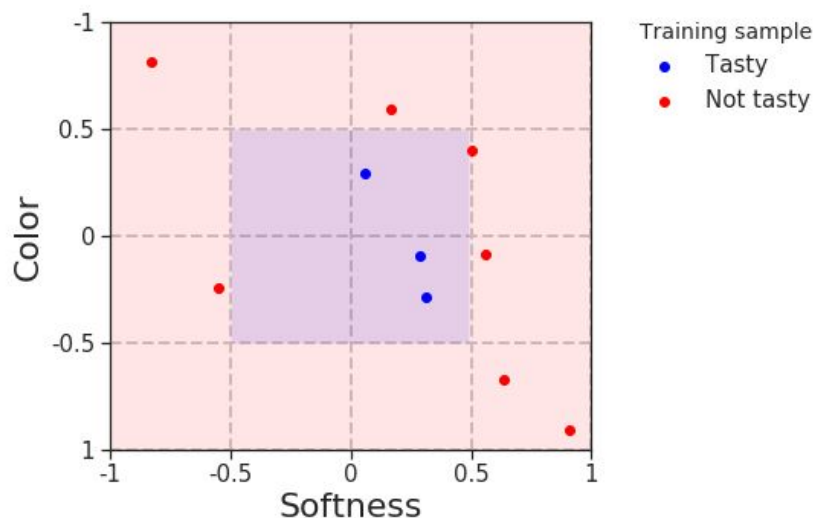
A controlled example:

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,

$x = [\text{softness}, \text{color}]$

Y : set of possible labels,

$y = [\text{tasty}, \text{not tasty}]$

D : data generation model,

$D \Rightarrow P(X)$

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

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

m : number of objects for training

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

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

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

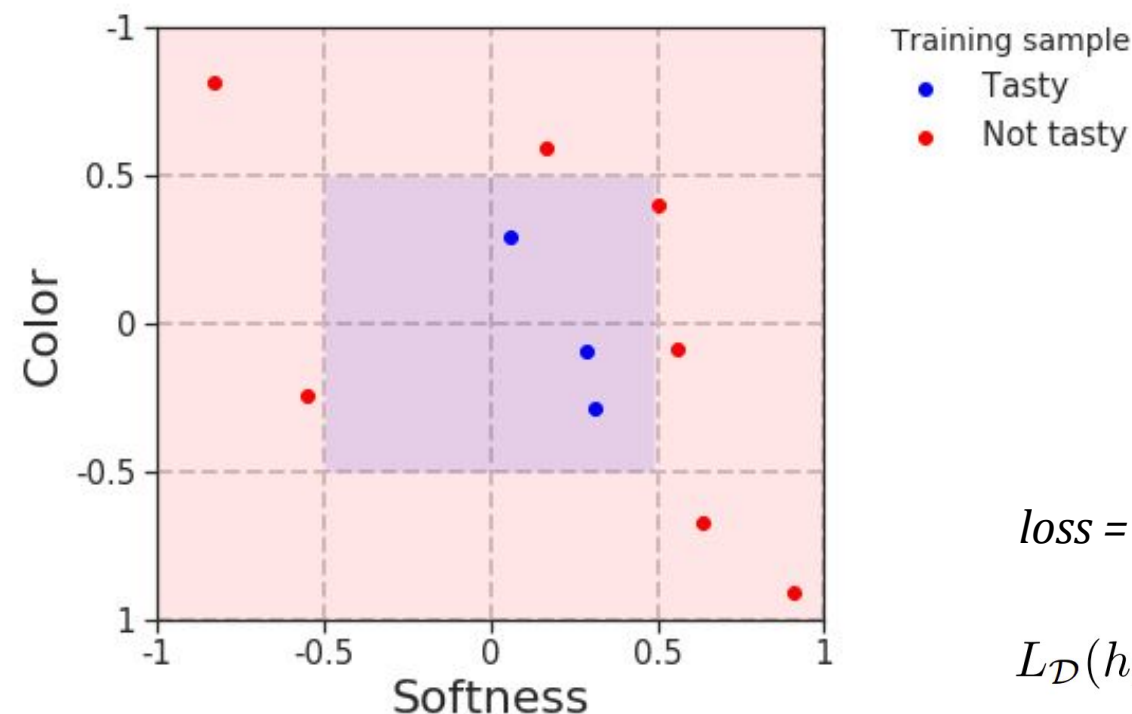
Question:

Proposed learner:

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

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

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



loss = fraction of incorrect predictions

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

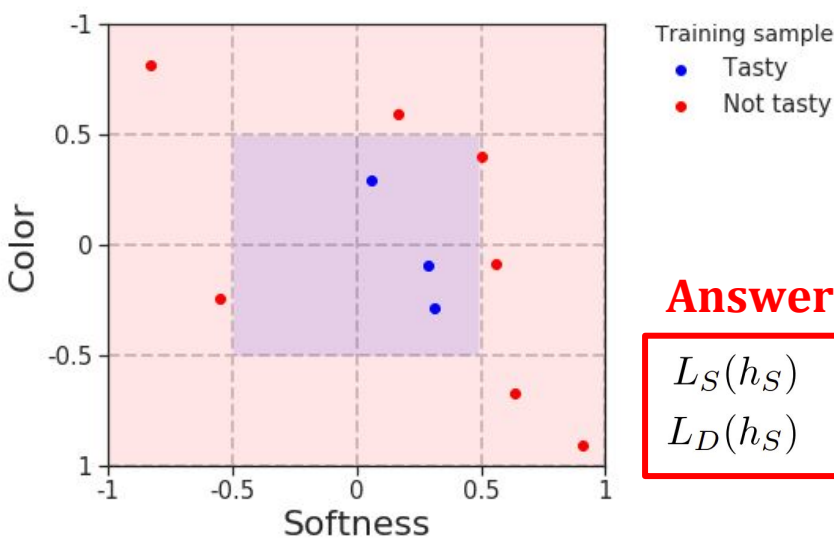
A controlled example:

Papaya tasting

Proposed learner:

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

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

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

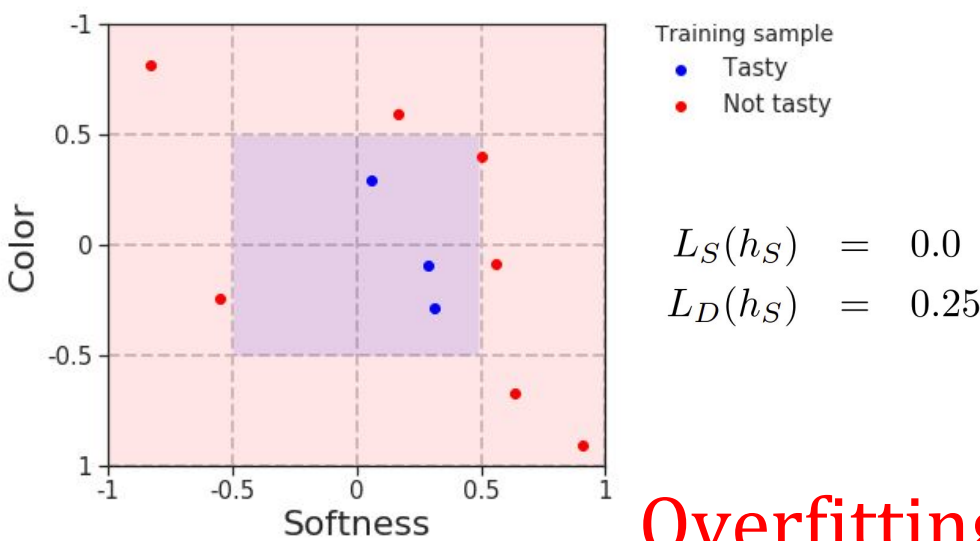
A controlled example:

Papaya tasting

Proposed learner:

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

Answer:



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

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

Overfitting! 🤔

Question:

- How can we avoid overfitting?

Question:

- How can we avoid overfitting?

by adding prior knowledge ...

Choosing the learner

X : set of all features,
 $x = [\text{softness}, \text{color}]$

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

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

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

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

m : number of objects for training

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

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

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

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

Hypothesis class (\mathcal{H}):

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

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

Choosing the learner

X : set of all features,
 $x = [\text{softness}, \text{color}]$

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

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

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

S : training sample: $[x_i, y_i], i \in \text{training}$
 h_S learner: $y_{\text{est};i} = h_S(x_i)$

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

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

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

Hypothesis class (\mathcal{H}):

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

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

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

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

ERM with inductive bias

\mathcal{X} : set of all features,

$x = [\text{softness}, \text{color}]$

\mathcal{Y} : set of possible labels,

$y = [\text{tasty}, \text{not tasty}]$

D : data generation model,

$D \Rightarrow P(\mathcal{X})$

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

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

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

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

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

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

Hypothesis class (\mathcal{H}):

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

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

- \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

- 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 S_1 is said to be representative of a probability distribution P if one can draw accurate conclusions about P from S_1
- If two samples S_1 and S_2 are representative of P , S_1 and S_2 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

\mathcal{X} : set of all features,

$x = [\text{softness}, \text{color}]$

\mathcal{Y} : set of possible labels,

$y = [\text{tasty}, \text{not tasty}]$

D : data generation model,

$D \Rightarrow P(\mathcal{X})$

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

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

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

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

Hypothesis class (\mathcal{H}):

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

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

- \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

- S is statistically representative of D

Things can still go wrong ...

Bad hypothesis and samples

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

$1 - \delta \rightarrow$ confidence parameter

Things can still go wrong ...

Bad hypothesis and samples

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

$1 - \delta \rightarrow$ confidence parameter

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

Good
hypothesis:

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

Bad
hypothesis:

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

Things can still go wrong ...

Bad hypothesis and samples

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

$1 - \delta \rightarrow$ confidence parameter

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

Good
hypothesis:

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

Bad
hypothesis:

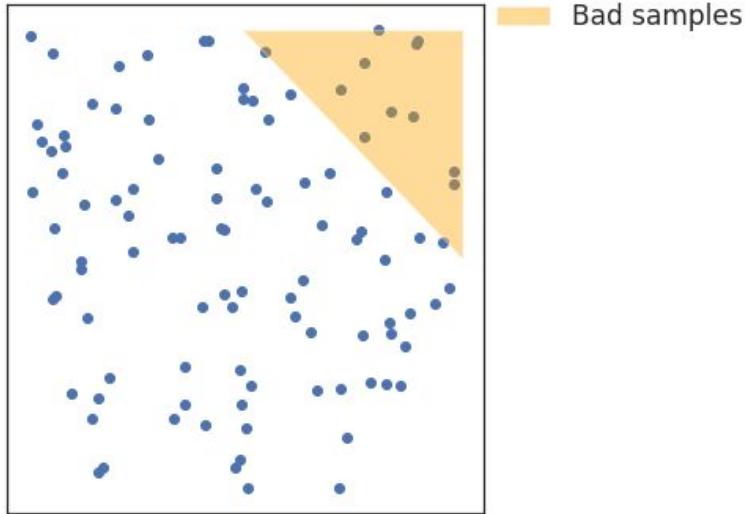
$$\mathcal{H}_B := [h \in \mathcal{H} : L_S(h_S) = 0 \quad \& \quad L_D(h_S) \geq \epsilon]$$

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

Things can still go wrong ...

Constructing misleading samples

The world



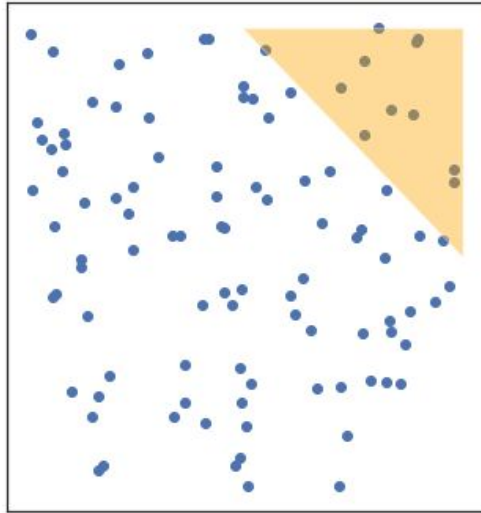
For 1 element in the training sample

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

Things can still go wrong ...

Constructing misleading samples

The world



Bad samples

For 1 element in the training sample

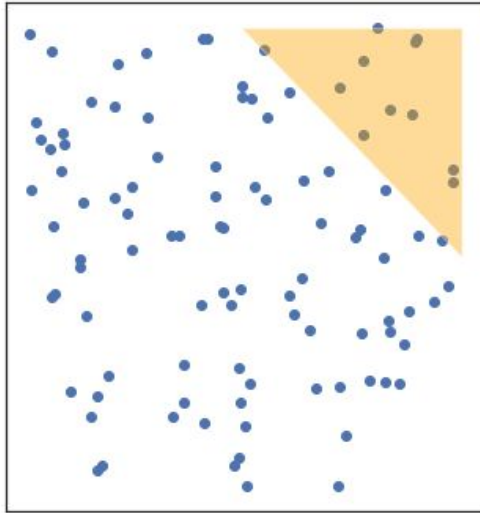
$$x_i \mid h(x_i) = y_i$$

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

Things can still go wrong ...

Constructing misleading samples

The world



Bad samples

For 1 element in the training sample

$$x_i \mid h(x_i) = y_i$$

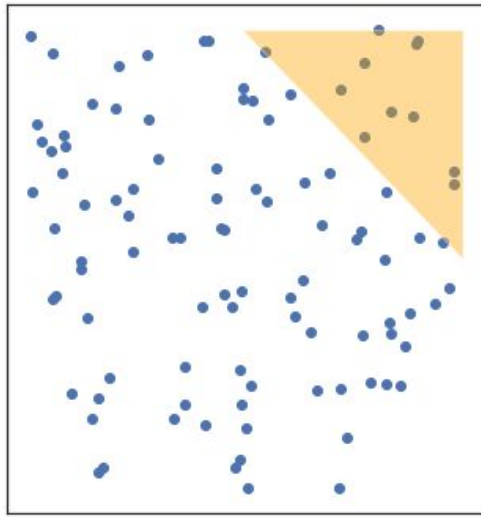
$$P(x_i \in \mathcal{D} : h(x_i) = y_i) = 1 - L_{\mathcal{D},f}(h)$$

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

Things can still go wrong ...

Constructing misleading samples

The world



Bad samples

For 1 element in the training sample

$$x_i \mid h(x_i) = y_i$$

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

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

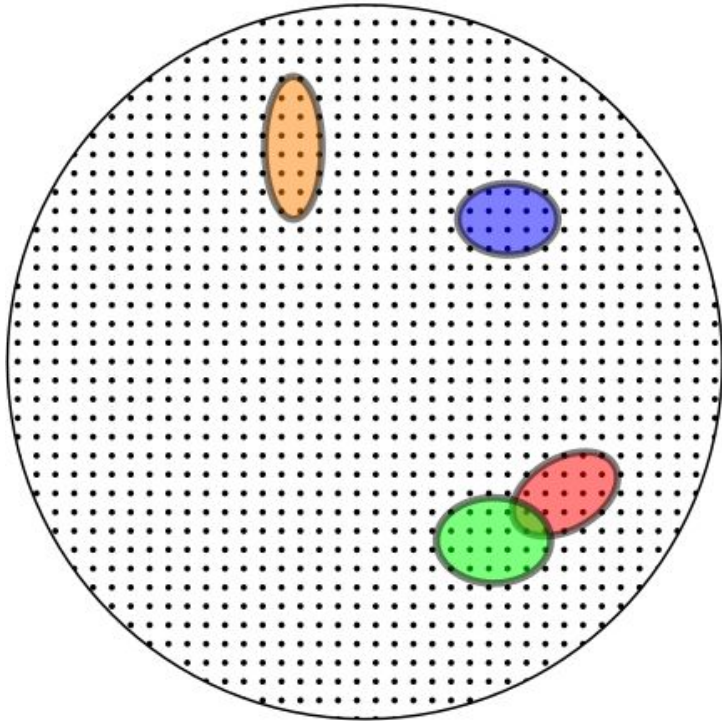
For m elements in the training sample

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

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

Things can still go wrong ...

Considering bad hypothesis

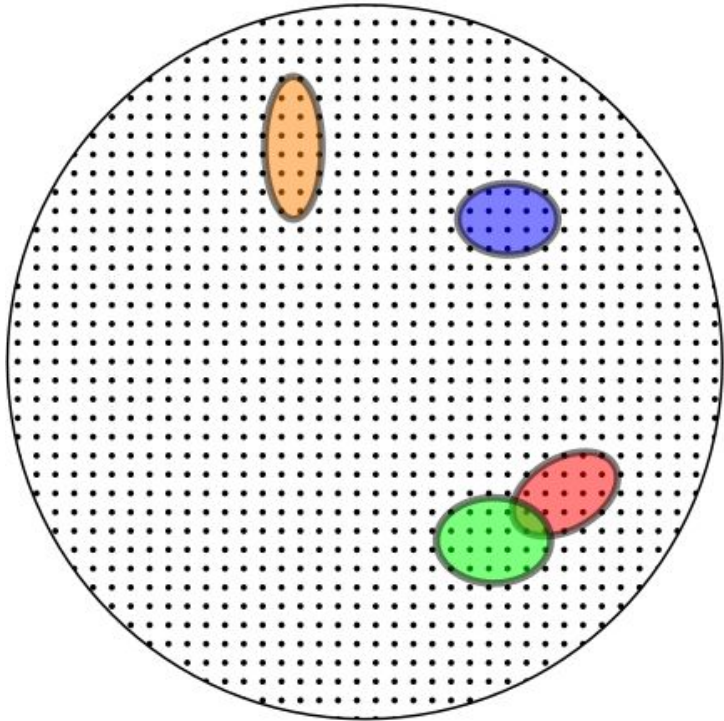


For 1 hypothesis

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

Things can still go wrong ...

Considering bad hypothesis



For 1 hypothesis

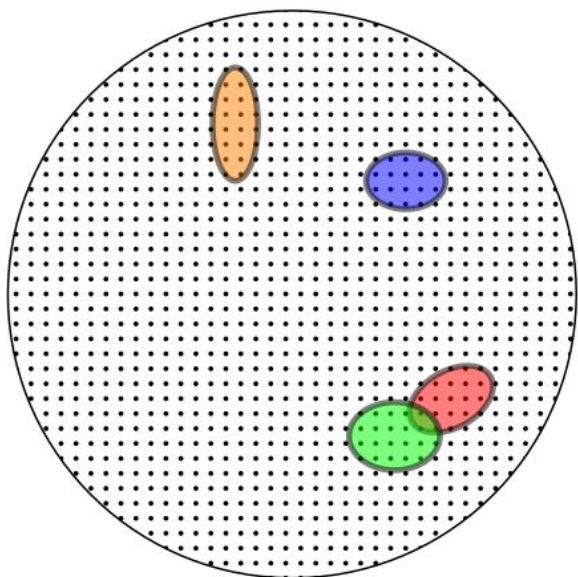
$$P(S_m : L_S(h) = 0) \leq (1 - \epsilon)^m$$

For all bad hypothesis

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

Things can still go wrong ...

Considering bad hypothesis



For 1 hypothesis

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

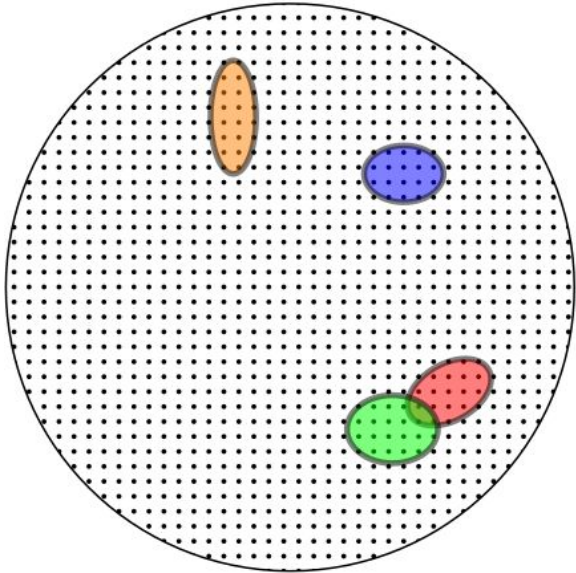
For all bad hypothesis

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

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

Things can still go wrong ...

Considering bad hypothesis



For 1 hypothesis

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

For all bad hypothesis

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

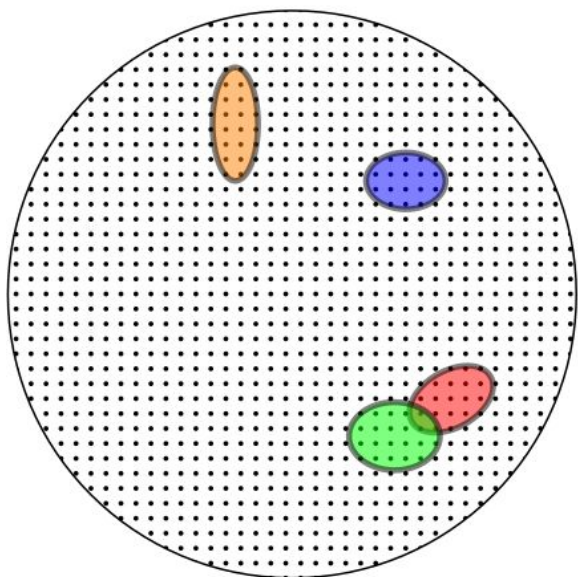
$$\delta = P(L_S(h) = 0, \forall h \in \mathcal{H}_B) \leq \sum_{h \in \mathcal{H}_B} (1 - \epsilon)^m$$

using...

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

Things can still go wrong ...

Considering bad hypothesis



For 1 hypothesis

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

For all bad hypothesis

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

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

using...

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

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

PAC learning model

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

Probably \rightarrow with confidence $1 - \delta$ over m samples
Approximately \rightarrow within a contamination level $\leq \epsilon$
Correct

If,

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

Return to a controlled example ...

Papaya tasting



χ : set of $x \in [\text{softness}, \text{color}]$

Y : set $y = [\text{tasty}, \text{not tasty}]$

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

Return to a controlled example ...

Papaya tasting



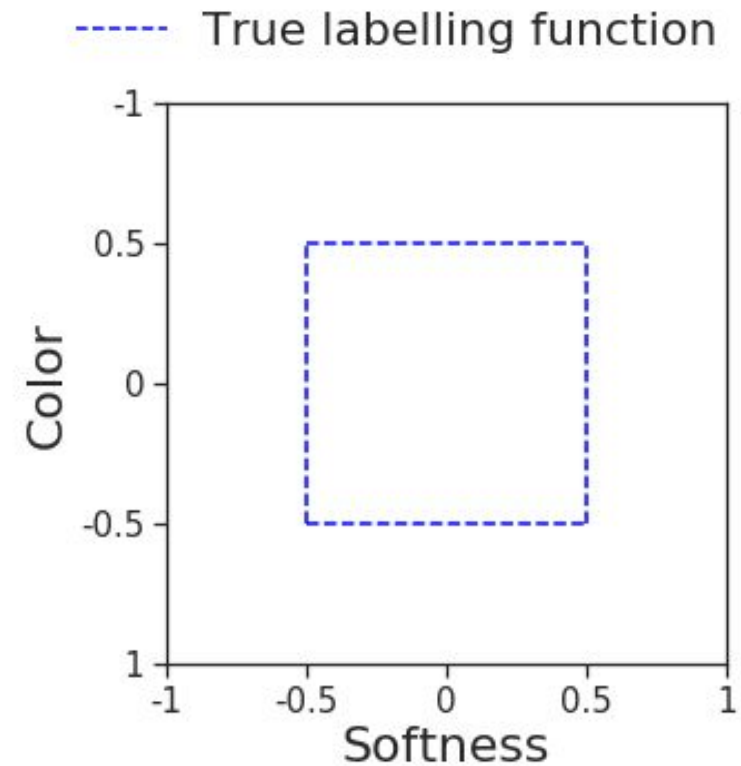
χ : set of $x \in [\text{softness}, \text{color}]$

Y : set $y = [\text{tasty}, \text{not tasty}]$

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

True Labelling function:

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



Return to a controlled example ...

Papaya tasting



χ : set of $x \in [\text{softness}, \text{color}]$

Y : set $y = [\text{tasty}, \text{not tasty}]$

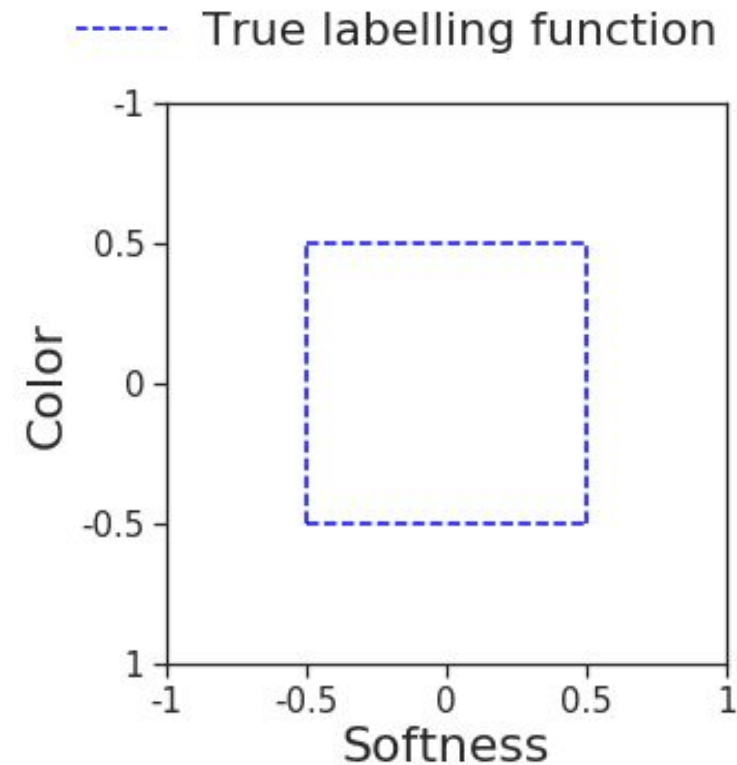
D : data generation model: $D \Rightarrow P(\chi)$

True Labelling function:

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

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

m : number of objects for training



Return to a controlled example ...

Papaya tasting



χ : set of $x \in [\text{softness}, \text{color}]$

Y : set $y = [\text{tasty}, \text{not tasty}]$

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

True Labelling function:

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

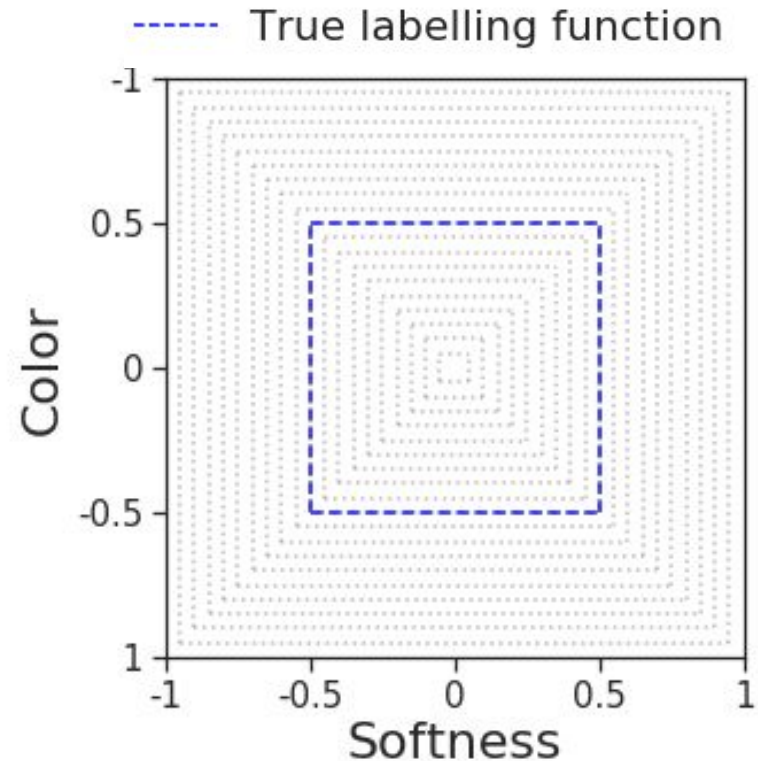
S : training sample: $[x_i, y_i], i \in \text{training}$

m : number of objects for training

\mathcal{H} : hypothesis class:

axis aligned squares in steps of 0.05

$N_H = 20$



Return to a controlled example ...

Papaya tasting



χ : set of $x \in [\text{softness}, \text{color}]$

Y : set $y = [\text{tasty}, \text{not tasty}]$

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

True Labelling function:

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

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

m : number of objects for training

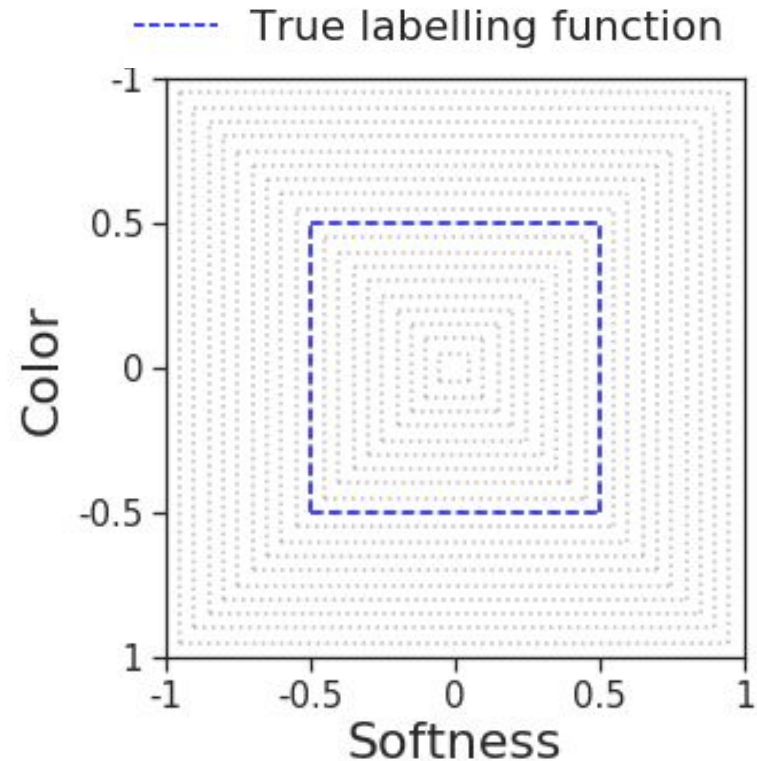
\mathcal{H} : hypothesis class:

axis aligned squares in steps of 0.05

$N_H = 20$

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

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



Return to a controlled example ...

Question:



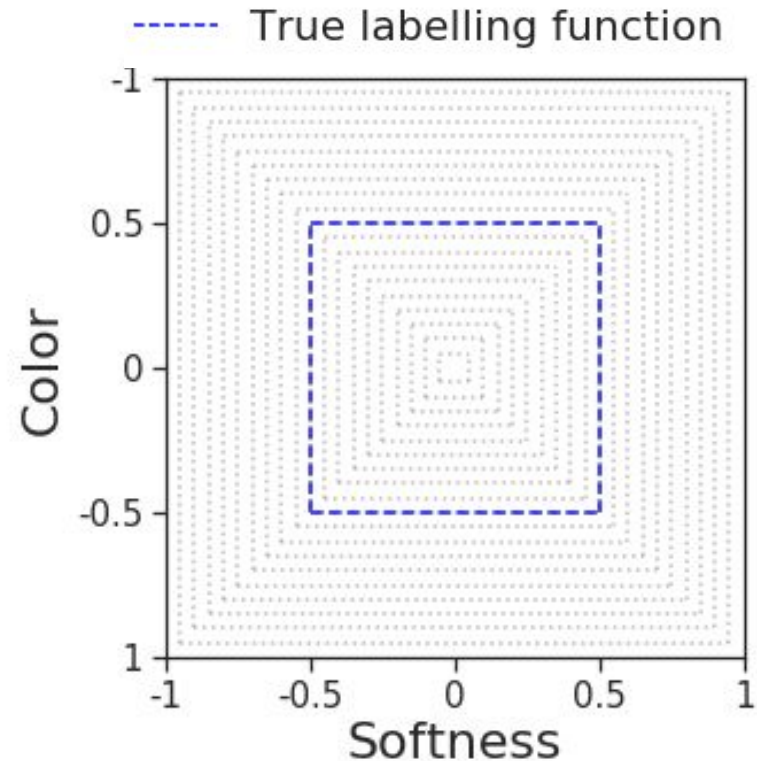
Data model: uniform distribution
[-1,1] in both axis

$1 - \delta = 0.95 \leftarrow$ confidence

$\epsilon = 0.05 \leftarrow$ contamination

$N_H = 20 \leftarrow$ number of possible
squares

m = ??



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

Return to a controlled example ...

Minimum number of samples

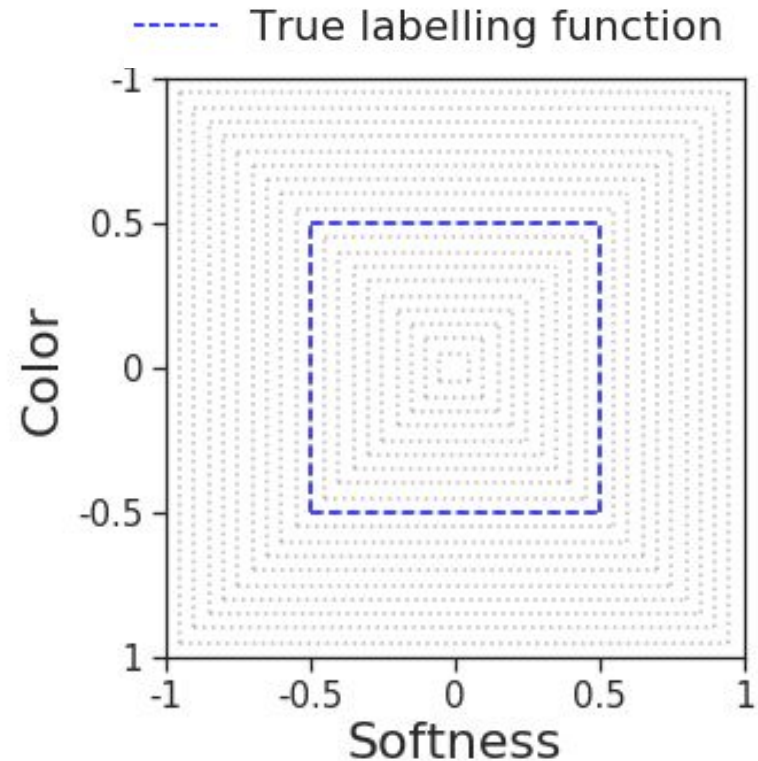
Data model: uniform distribution
[-1,1] in both axis

$1 - \delta = 0.95 \leftarrow$ confidence

$\varepsilon = 0.05 \leftarrow$ contamination

$N_H = 20 \leftarrow$ number of possible
squares

$m \sim 120$

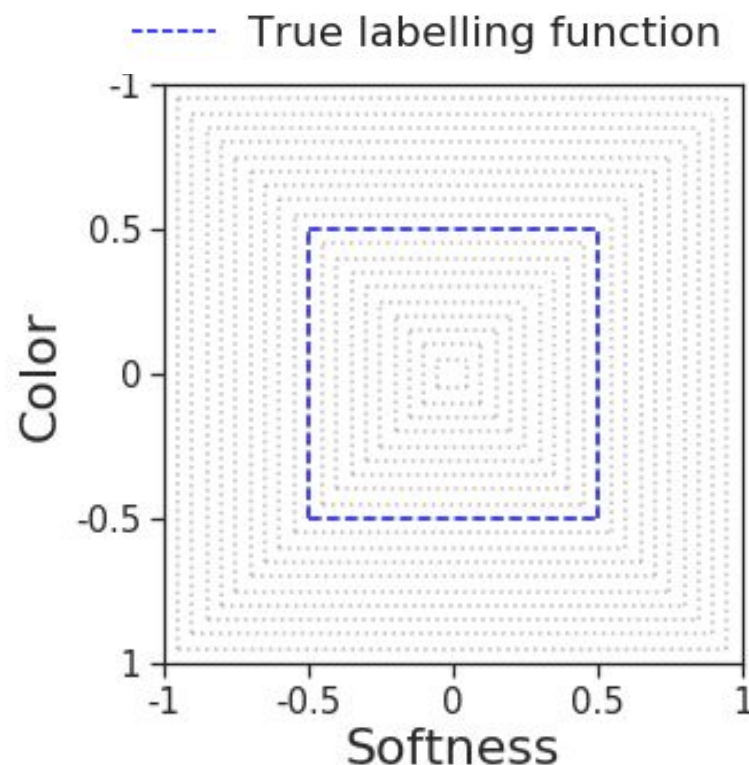
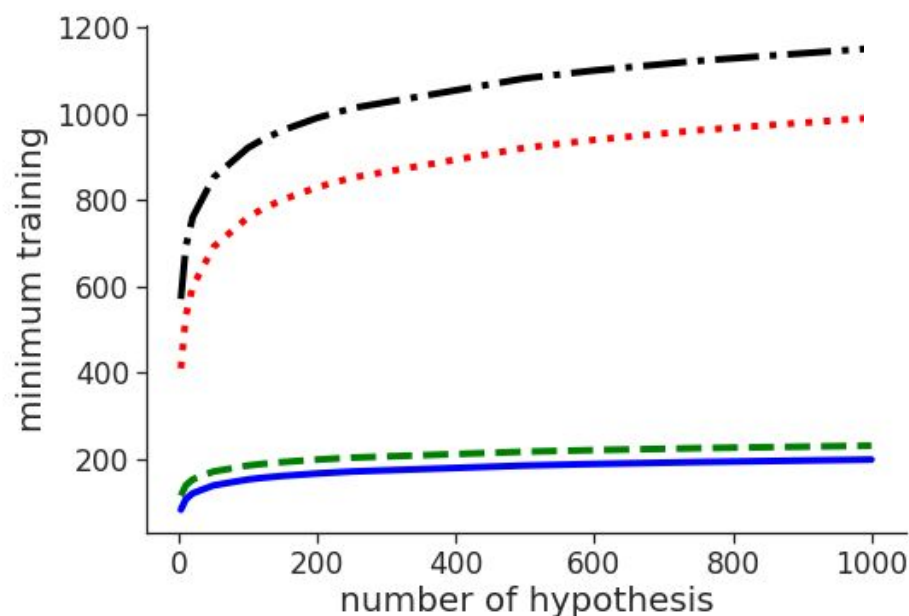


Return to a controlled example ...

Minimum number of samples

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

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



— · — $\delta = 0.01, \varepsilon = 0.01$
... $\delta = 0.05, \varepsilon = 0.01$
- - - $\delta = 0.01, \varepsilon = 0.05$
— $\delta = 0.05, \varepsilon = 0.05$

If you still have stomach...

Next theory session:

- Agnostic PAC learning
- Uniform convergence
- Infinity number of hypothesis ...

Thank you!