

PAC Learning

CPT_S 434/534 Neural network design and application

Core questions to answer


- What can be learned by machine learning models?
- What conditions are required to successfully learn?

Underlying concept function

concept $c: \mathcal{X} \rightarrow \mathcal{Y}$

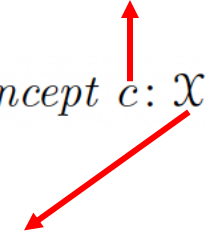
Underlying concept function

An underlying concept function (mapping)

$$\text{concept } c: \mathcal{X} \rightarrow \mathcal{Y}$$


Underlying concept function

An underlying concept function (mapping)

$$\text{concept } c: \mathcal{X} \rightarrow \mathcal{Y}$$
A diagram consisting of two red arrows. One arrow points vertically upwards from the text 'concept c: X -> Y' to the text 'An underlying concept function (mapping)'. The other arrow points diagonally downwards and to the left from the text 'concept c: X -> Y' to the text 'Some feature space'.

Some feature space

Underlying concept function

An underlying concept function (mapping)

$$\text{concept } c: \mathcal{X} \rightarrow \mathcal{Y}$$

Some feature space

$$\mathcal{Y} = \{0, 1\}$$

$$\text{or } \mathcal{Y} = \mathbb{R}$$

Underlying concept function

An underlying concept function (mapping)

$$\text{concept } c: \mathcal{X} \rightarrow \mathcal{Y}$$

Some feature space

$$\mathcal{Y} = \{0, 1\}$$

or $\mathcal{Y} = \mathbb{R}$



Features

Size
bedrooms
Yard
.....

A model

Price?

Underlying concept function

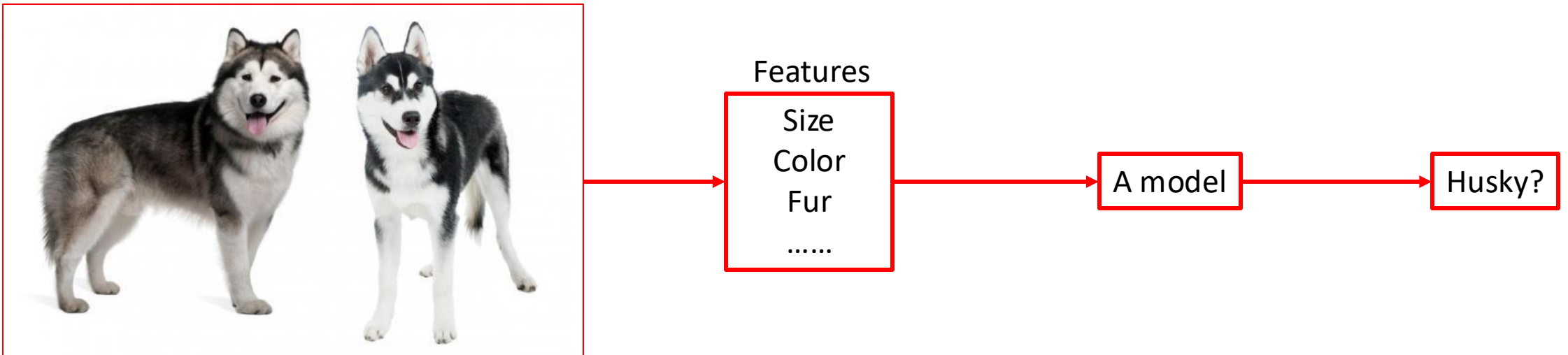
An underlying concept function (mapping)

$$\text{concept } c: \mathcal{X} \rightarrow \mathcal{Y}$$

Some feature space

$$\mathcal{Y} = \{0, 1\}$$

or $\mathcal{Y} = \mathbb{R}$



Underlying concept function

An underlying concept function (mapping)

$$\text{concept } c: \mathcal{X} \rightarrow \mathcal{Y}$$

Some feature space

$$\mathcal{Y} = \{0, 1\}$$

or $\mathcal{Y} = \mathbb{R}$



Features

Color
Shape
Yellow spot
.....

A model

Sweet?

Underlying concept function

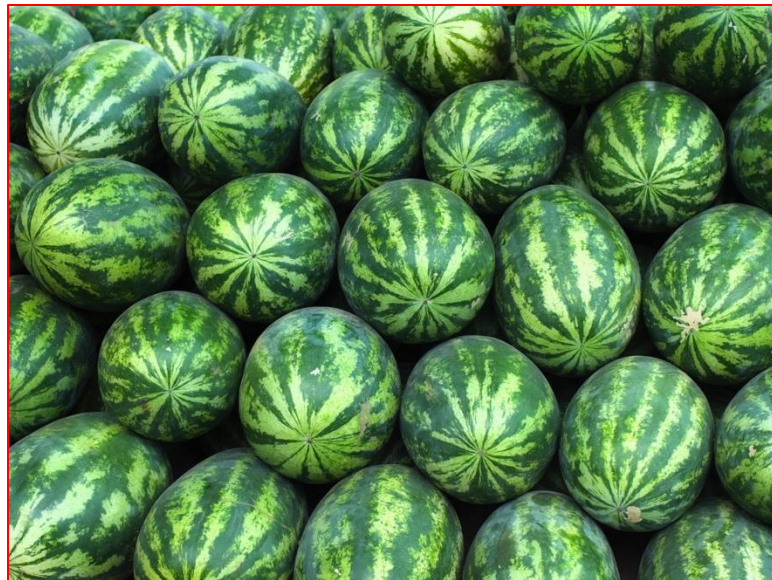
An underlying concept function (mapping)

$$\text{concept } c: \mathcal{X} \rightarrow \mathcal{Y}$$

Some feature space

$$\mathcal{Y} = \{0, 1\}$$

or $\mathcal{Y} = \mathbb{R}$



Features

Color
Shape
Yellow spot
.....

A model

Sweet?

Q: how to measure the model's performance?

Underlying concept function

An underlying concept function (mapping)

$$\text{concept } c: \mathcal{X} \rightarrow \mathcal{Y}$$

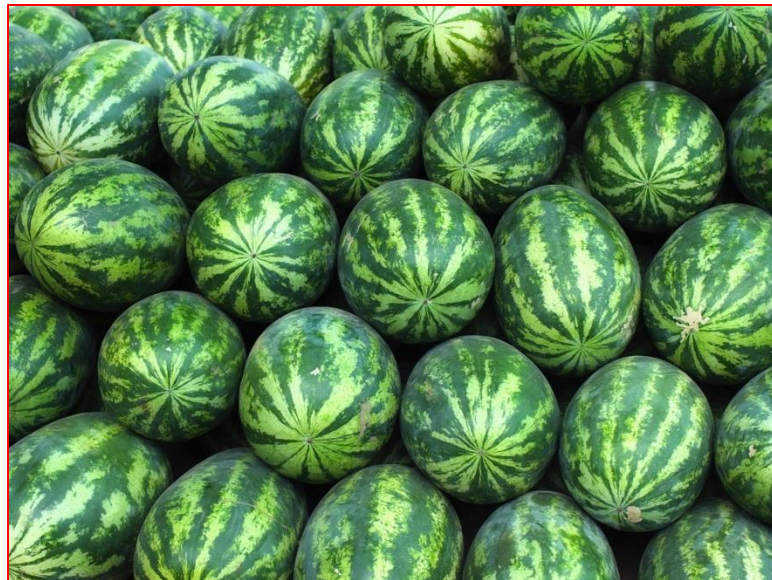
Some feature space

$$\mathcal{Y} = \{0, 1\}$$

or $\mathcal{Y} = \mathbb{R}$

Q: What if our model has **completely different** behavior with the underlying concept function?

Q: how to measure the model's performance?



Features

Color
Shape
Yellow spot
.....

A model

Sweet?

Underlying concept function

An underlying concept function (mapping)

$$\text{concept } c: \mathcal{X} \rightarrow \mathcal{Y}$$

Some feature space

$$\mathcal{Y} = \{0, 1\}$$

or $\mathcal{Y} = \mathbb{R}$

Q: What if our model has **identical** behavior with the underlying concept function?

Q: how to measure the model's performance?



Features

Color
Shape
Yellow spot
.....

A model

Sweet?

Risk

A hypothesis class

Definition 2.1 (Generalization error) *Given a hypothesis $h \in \mathcal{H}$, a target concept $c \in \mathcal{C}$, and an underlying distribution \mathcal{D} , the generalization error or risk of h is defined by*

$$R(h) = \mathbb{P}_{x \sim \mathcal{D}} [h(x) \neq c(x)] = \mathbb{E}_{x \sim \mathcal{D}} [1_{h(x) \neq c(x)}], \quad (2.1)$$

where 1_ω is the indicator function of the event ω .²

Risk

A hypothesis class



Definition 2.1 (Generalization error) *Given a hypothesis $h \in \mathcal{H}$, a target concept $c \in \mathcal{C}$, and an underlying distribution \mathcal{D} , the generalization error or risk of h is defined by*

$$R(h) = \mathbb{P}_{x \sim \mathcal{D}} [h(x) \neq c(x)] = \mathbb{E}_{x \sim \mathcal{D}} [1_{h(x) \neq c(x)}], \quad (2.1)$$

where 1_ω is the indicator

Risk: in population level
→
scan all samples in the world
(not feasible in general)

Risk

Q: What is it?

A hypothesis class



Definition 2.1 (Generalization error) *Given a hypothesis $h \in \mathcal{H}$, a target concept $c \in \mathcal{C}$, and an underlying distribution \mathcal{D} , the generalization error or risk of h is defined by*

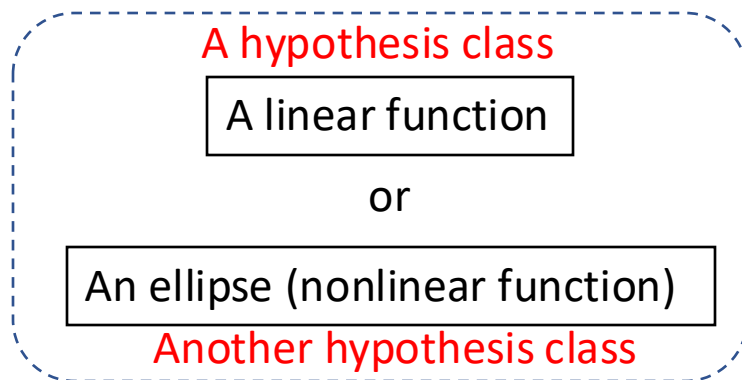
$$R(h) = \mathbb{P}_{x \sim \mathcal{D}} [h(x) \neq c(x)] = \mathbb{E}_{x \sim \mathcal{D}} [1_{h(x) \neq c(x)}], \quad (2.1)$$

where 1_ω is the indicator

Risk: in population level
→
scan all samples in the world
(not feasible in general)

Review: Build a model

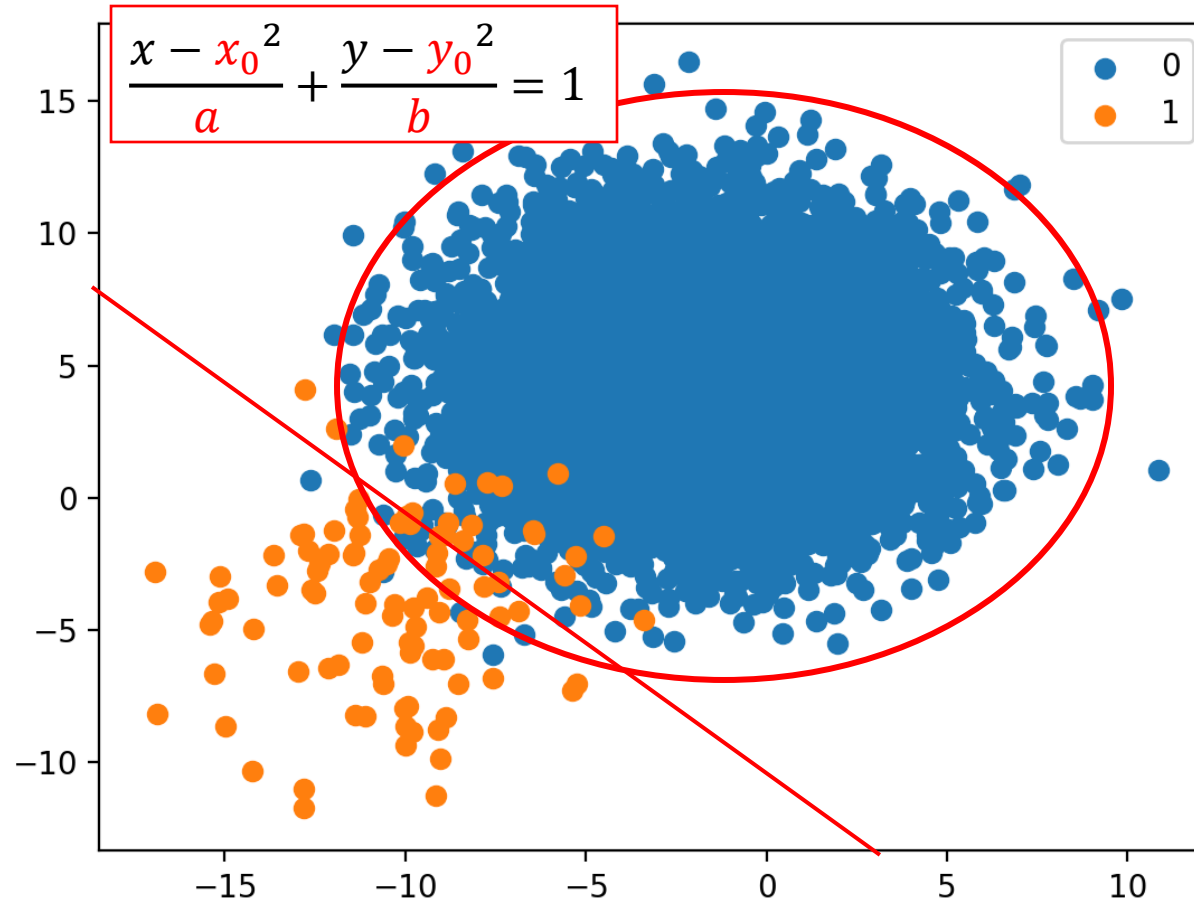
- What is a model



Q: what are their parameters?

$$y = ax + b$$

Try to separate two classes
Q: how to separate them?



Risk

Q: What is it?

A hypothesis class

Definition 2.1 (Generalization error) Given a hypothesis $h \in \mathcal{H}$, a target concept $c \in \mathcal{C}$, and an underlying distribution \mathcal{D} , the generalization error or risk of h is defined by

$$R(h) = \mathbb{P}_{x \sim \mathcal{D}} [h(x) \neq c(x)] = \mathbb{E}_{x \sim \mathcal{D}} [1_{h(x) \neq c(x)}], \quad (2.1)$$

where 1_ω is the indicator

Risk: in population level
→
scan all samples in the world
(not feasible in general)

Expected mistakes that h
makes over data distribution D

PAC learning

Definition 2.3 (PAC-learning) *A concept class \mathcal{C} is said to be PAC-learnable if there exists an algorithm \mathcal{A} and a polynomial function $\text{poly}(\cdot, \cdot, \cdot, \cdot)$ such that for any $\epsilon > 0$ and $\delta > 0$ for all distributions \mathcal{D} on \mathcal{X} and for any target concept $c \in \mathcal{C}$, the following holds for any sample size $m \geq \text{poly}(1/\epsilon, 1/\delta, n, \text{size}(c))$:*

$$\mathbb{P}_{S \sim \mathcal{D}^m} [R(h_S) \leq \epsilon] \geq 1 - \delta. \quad (2.4)$$

If \mathcal{A} further runs in $\text{poly}(1/\epsilon, 1/\delta, n, \text{size}(c))$, then \mathcal{C} is said to be efficiently PAC-learnable. When such an algorithm \mathcal{A} exists, it is called a PAC-learning algorithm for \mathcal{C} .

PAC learning

Definition 2.3 (PAC-learning) A concept class \mathcal{C} is said to be PAC-learnable if there exists an algorithm \mathcal{A} and a polynomial function $\text{poly}(\cdot, \cdot, \cdot, \cdot)$ such that for any $\epsilon > 0$ and $\delta > 0$ for all distributions \mathcal{D} on \mathcal{X} and for any target concept $c \in \mathcal{C}$, the following holds for any sample size $m \geq \text{poly}(1/\epsilon, 1/\delta, n, \text{size}(c))$:

$$\mathbb{P}_{S \sim \mathcal{D}^m} [R(h_S) \leq \epsilon] \geq 1 - \delta. \quad (2.4)$$

If \mathcal{A} further runs in $\text{poly}(1/\epsilon, 1/\delta, n, \text{size}(c))$, then \mathcal{C} is said to be efficiently PAC-learnable. When such an algorithm \mathcal{A} exists, it is called a PAC-learning algorithm for \mathcal{C} .

PAC learning

Definition 2.3 (PAC-learning) A concept class \mathcal{C} is said to be PAC-learnable if there exists an algorithm \mathcal{A} and a polynomial function $\text{poly}(\cdot, \cdot, \cdot, \cdot)$ such that for any $\epsilon > 0$ and $\delta > 0$, for all distributions \mathcal{D} on \mathcal{X} and for any target concept $c \in \mathcal{C}$, the following holds for any sample size $m \geq \text{poly}(1/\epsilon, 1/\delta, n, \text{size}(c))$:

$$\mathbb{P}_{S \sim \mathcal{D}^m} [R(h_S) \leq \epsilon] \geq 1 - \delta. \quad \text{Probably} \quad (2.4)$$

Approximately **C**orrect

If \mathcal{A} further runs in $\text{poly}(1/\epsilon, 1/\delta, n, \text{size}(c))$, then \mathcal{C} is said to be efficiently PAC-learnable. When such an algorithm \mathcal{A} exists, it is called a PAC-learning algorithm for \mathcal{C} .

PAC learning

Polynomial: $a_n x^n + a_{n-1} x^{n-1} + \dots + a_2 x^2 + a_1 x + a_0,$

Definition 2.3 (PAC-learning) *A concept class \mathcal{C} is said to be PAC-learnable if there exists an algorithm \mathcal{A} and a polynomial function $\text{poly}(\cdot, \cdot, \cdot, \cdot)$ such that for any $\epsilon > 0$ and $\delta > 0$, for all distributions \mathcal{D} on \mathcal{X} and for any target concept $c \in \mathcal{C}$, the following holds for any sample size $m \geq \text{poly}(1/\epsilon, 1/\delta, n, \text{size}(c))$:*

$$\mathbb{P}_{S \sim \mathcal{D}^m} [R(h_S) \leq \epsilon] \geq 1 - \delta. \quad (2.4)$$

If \mathcal{A} further runs in $\text{poly}(1/\epsilon, 1/\delta, n, \text{size}(c))$, then \mathcal{C} is said to be efficiently PAC-learnable. When such an algorithm \mathcal{A} exists, it is called a PAC-learning algorithm for \mathcal{C} .

PAC learning

Polynomial: $a_n x^n + a_{n-1} x^{n-1} + \dots + a_2 x^2 + a_1 x + a_0,$

Definition 2.3 (PAC-learning) A concept class \mathcal{C} is said to be PAC-learnable if there exists an algorithm \mathcal{A} and a polynomial function $\text{poly}(\cdot, \cdot, \cdot, \cdot)$ such that for any $\epsilon > 0$ and $\delta > 0$, for all distributions \mathcal{D} on \mathcal{X} and for any target concept $c \in \mathcal{C}$, the following holds for any sample size $m \geq \text{poly}(1/\epsilon, 1/\delta, n, \text{size}(c))$:

$$\mathbb{P}_{S \sim \mathcal{D}^m} [R(h_S) \leq \epsilon] \geq 1 - \delta. \quad (2.4)$$

If \mathcal{A} further runs in $\text{poly}(1/\epsilon, 1/\delta, n, \text{size}(c))$, then \mathcal{C} is said to be efficiently PAC-learnable. When such an algorithm \mathcal{A} exists, it is called a PAC-learning algorithm for \mathcal{C} .

$$m \rightarrow S \rightarrow h_S$$

Underlying concept function

An underlying concept function (mapping)

$$\text{concept } c: \mathcal{X} \rightarrow \mathcal{Y}$$

Some feature space

$$\mathcal{Y} = \{0, 1\}$$

or $\mathcal{Y} = \mathbb{R}$



Features

Color
Shape
Yellow spot
.....

A model

Sweet?

A new Q: with almost the same features, can we guarantee the prediction (sweetness)?

Q: how to measure the model's performance?

Underlying concept function

An underlying concept function (mapping)

$$\text{concept } c: \mathcal{X} \rightarrow \mathcal{Y}$$

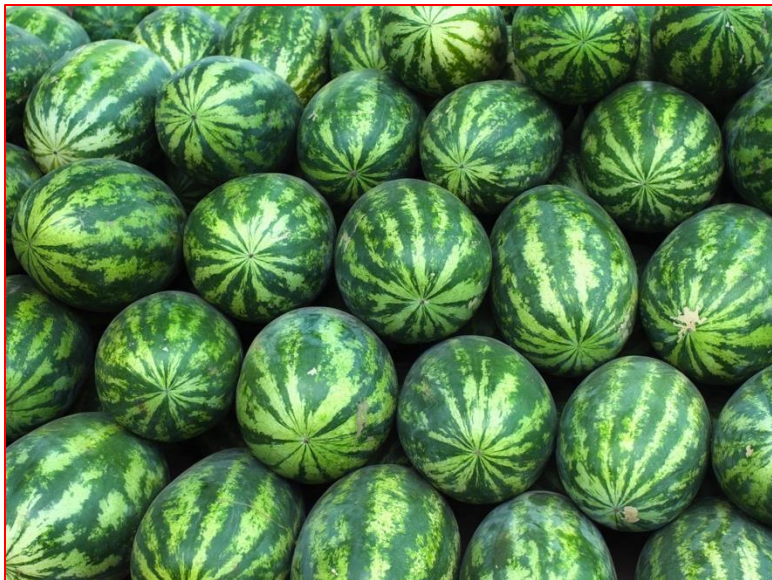
Some feature space

$$\mathcal{Y} = \{0, 1\}$$

or $\mathcal{Y} = \mathbb{R}$

Nearly the same color/shape/...
→
Different sweetness?

A new Q: with almost the same features, can we guarantee the prediction (sweetness)?



Features

Color
Shape
Yellow spot
.....

A model

Sweet?

Q: how to measure the model's performance?

Agnostic PAC learning

Definition 2.14 (Agnostic PAC-learning) Let \mathcal{H} be a hypothesis set. \mathcal{A} is an agnostic PAC-learning algorithm if there exists a polynomial function $\text{poly}(\cdot, \cdot, \cdot, \cdot)$ such that for any $\epsilon > 0$ and $\delta > 0$, for all distributions \mathcal{D} over $\mathcal{X} \times \mathcal{Y}$, the following holds for any sample size $m \geq \text{poly}(1/\epsilon, 1/\delta, n, \text{size}(c))$:

$$\mathbb{P}_{S \sim \mathcal{D}^m} [R(h_S) - \min_{h \in \mathcal{H}} R(h) \leq \epsilon] \geq 1 - \delta. \quad (2.21)$$

If \mathcal{A} further runs in $\text{poly}(1/\epsilon, 1/\delta, n)$, then it is said to be an efficient agnostic PAC-learning algorithm.

Agnostic PAC learning

$$R(h) = \mathbb{P}_{x \sim \mathcal{D}} [h(x) \neq c(x)] = \mathbb{E}_{x \sim \mathcal{D}} [1_{h(x) \neq c(x)}],$$

$c(x)$ is deterministic

$$R(h) = \mathbb{P}_{(x,y) \sim \mathcal{D}} [h(x) \neq y] = \mathbb{E}_{(x,y) \sim \mathcal{D}} [1_{h(x) \neq y}].$$

stochastic: joint distribution \mathcal{D}

Definition 2.14 (Agnostic PAC-learning) Let \mathcal{H} be a hypothesis set. \mathcal{A} is an agnostic PAC-learning algorithm if there exists a polynomial function $\text{poly}(\cdot, \cdot, \cdot, \cdot)$ such that for any $\epsilon > 0$ and $\delta > 0$, for all distributions \mathcal{D} over $\mathcal{X} \times \mathcal{Y}$, the following holds for any sample size $m \geq \text{poly}(1/\epsilon, 1/\delta, n, \text{size}(c))$:

$$\mathbb{P}_{S \sim \mathcal{D}^m} [R(h_S) - \min_{h \in \mathcal{H}} R(h) \leq \epsilon] \geq 1 - \delta. \quad (2.21)$$

If \mathcal{A} further runs in $\text{poly}(1/\epsilon, 1/\delta, n)$, then it is said to be an efficient agnostic PAC-learning algorithm.

Agnostic PAC learning

$$R(h) = \mathbb{P}_{x \sim \mathcal{D}} [h(x) \neq c(x)] = \mathbb{E}_{x \sim \mathcal{D}} [1_{h(x) \neq c(x)}],$$

$c(x)$ is deterministic

$$R(h) = \mathbb{P}_{(x,y) \sim \mathcal{D}} [h(x) \neq y] = \mathbb{E}_{(x,y) \sim \mathcal{D}} [1_{h(x) \neq y}].$$

stochastic: joint distribution \mathcal{D}

Definition 2.14 (Agnostic PAC-learning) Let \mathcal{H} be a hypothesis set. \mathcal{A} is an agnostic PAC-learning algorithm if there exists a polynomial function $\text{poly}(\cdot, \cdot, \cdot, \cdot)$ such that for any $\epsilon > 0$ and $\delta > 0$, for all distributions \mathcal{D} over $\mathcal{X} \times \mathcal{Y}$, the following holds for any sample size $m \geq \text{poly}(1/\epsilon, 1/\delta, n, \text{size}(c))$:

$$\mathbb{P}_{S \sim \mathcal{D}^m} [R(h_S) - \min_{h \in \mathcal{H}} R(h) \leq \epsilon] \geq 1 - \delta. \quad (2.21)$$

If \mathcal{A} further runs in $\text{poly}(1/\epsilon, 1/\delta, n)$, then it is said to be an efficient agnostic PAC-learning algorithm.

Bayes error

Definition 2.15 (Bayes error) Given a distribution \mathcal{D} over $\mathcal{X} \times \mathcal{Y}$, the Bayes error R^* is defined as the infimum of the errors achieved by measurable functions $h: \mathcal{X} \rightarrow \mathcal{Y}$:

$$R^* = \inf_{\substack{h \\ h \text{ measurable}}} R(h). \quad (2.22)$$

A hypothesis h with $R(h) = R^*$ is called a Bayes hypothesis or Bayes classifier.

All possible hypotheses
(may not be included in H)

Bayes error

Definition 2.15 (Bayes error) *Given a distribution \mathcal{D} over $\mathcal{X} \times \mathcal{Y}$, the Bayes error R^* is defined as the infimum of the errors achieved by measurable functions $h: \mathcal{X} \rightarrow \mathcal{Y}$:*

The **best** risk we may reach $\leftarrow R^* = \inf_{h \text{ measurable}} R(h). \quad (2.22)$

A hypothesis h with $R(h) = R^$ is called a Bayes hypothesis or Bayes classifier.*

All possible hypotheses
(may not be included in H)

Estimation and approximation

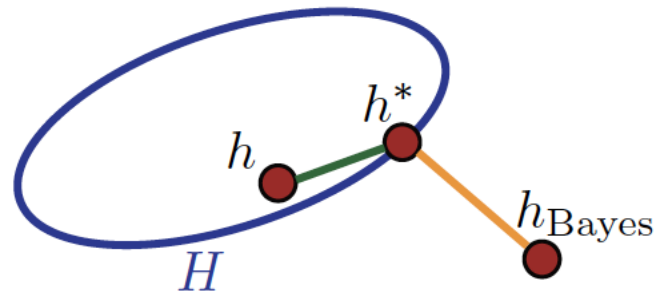
$$R(h) - R^*$$

Estimation and approximation

$$R(h) - R^* = \underbrace{\left(R(h) - \inf_{h \in \mathcal{H}} R(h) \right)}_{\text{estimation}} + \underbrace{\left(\inf_{h \in \mathcal{H}} R(h) - R^* \right)}_{\text{approximation}}.$$

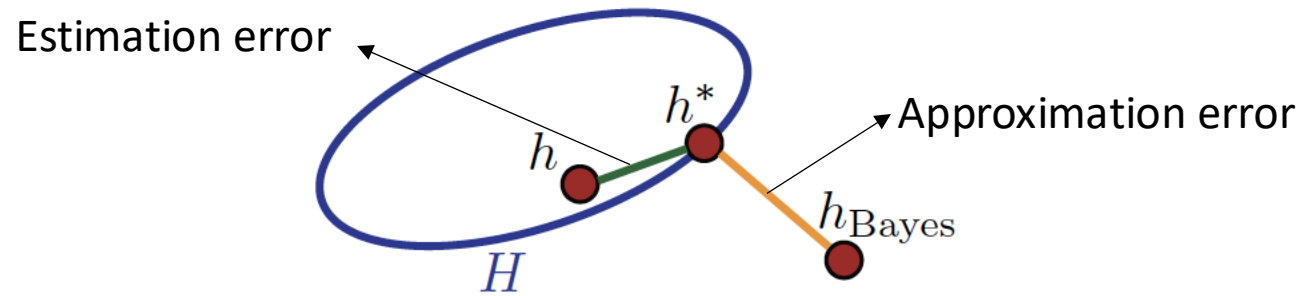
Estimation and approximation

$$R(h) - R^* = \underbrace{\left(R(h) - \inf_{h \in \mathcal{H}} R(h) \right)}_{\text{estimation}} + \underbrace{\left(\inf_{h \in \mathcal{H}} R(h) - R^* \right)}_{\text{approximation}}.$$



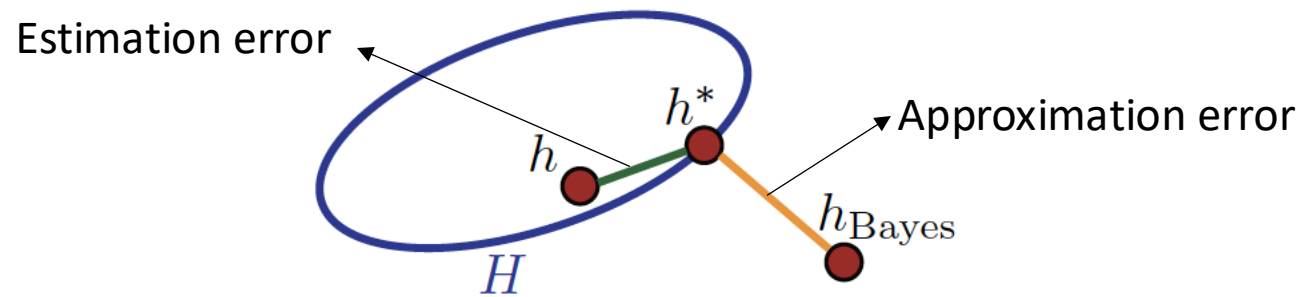
Estimation and approximation

$$R(h) - R^* = \underbrace{\left(R(h) - \inf_{h \in \mathcal{H}} R(h) \right)}_{\text{estimation}} + \underbrace{\left(\inf_{h \in \mathcal{H}} R(h) - R^* \right)}_{\text{approximation}}.$$



Estimation and approximation

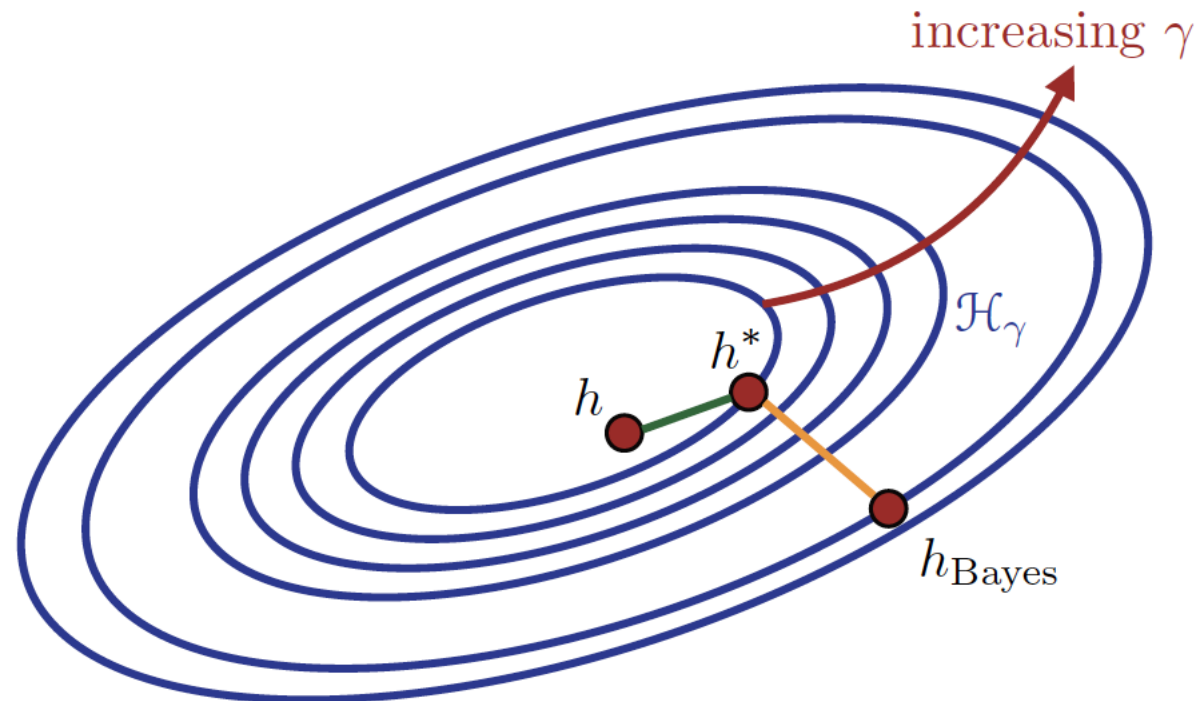
$$R(h) - R^* = \underbrace{\left(R(h) - \inf_{h \in \mathcal{H}} R(h) \right)}_{\text{estimation}} + \underbrace{\left(\inf_{h \in \mathcal{H}} R(h) - R^* \right)}_{\text{approximation}}.$$



Q: can we enlarge H ?

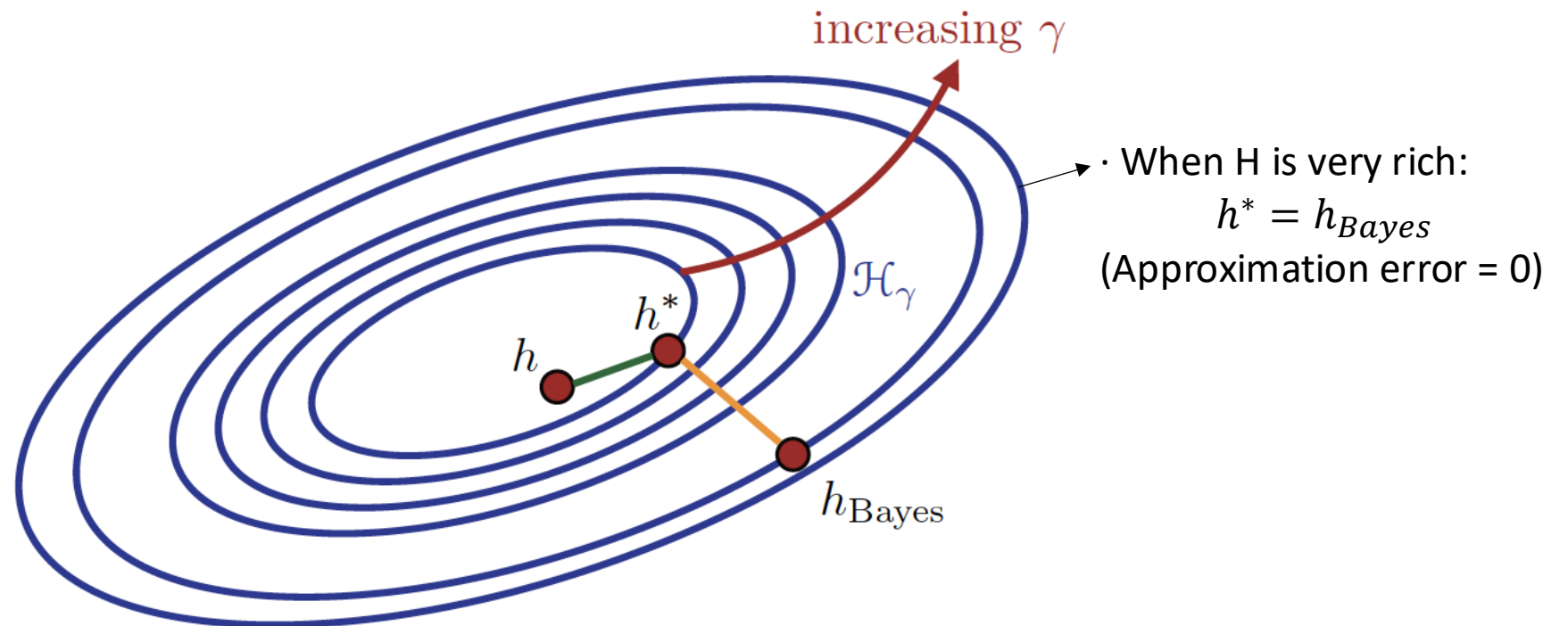
Estimation and approximation

$$R(h) - R^* = \underbrace{\left(R(h) - \inf_{h \in \mathcal{H}} R(h) \right)}_{\text{estimation}} + \underbrace{\left(\inf_{h \in \mathcal{H}} R(h) - R^* \right)}_{\text{approximation}}.$$



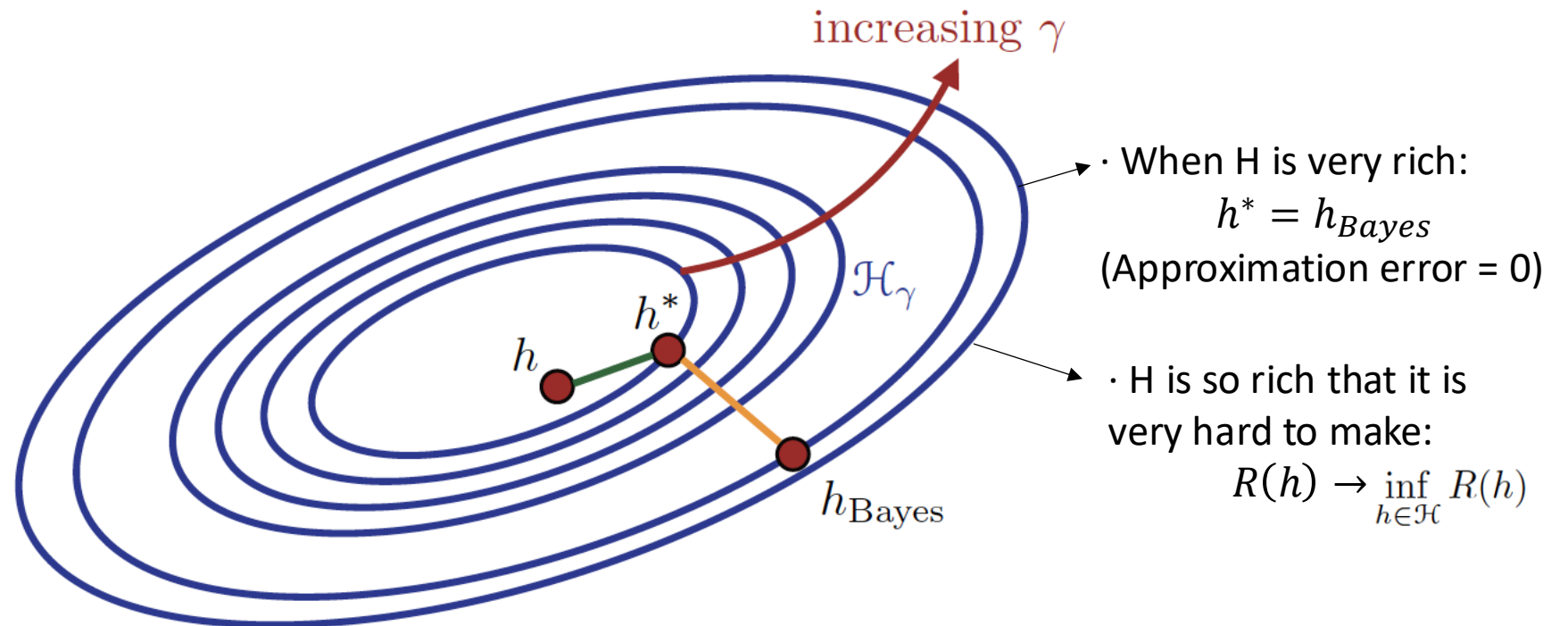
Estimation and approximation

$$R(h) - R^* = \underbrace{\left(R(h) - \inf_{h \in \mathcal{H}} R(h) \right)}_{\text{estimation}} + \underbrace{\left(\inf_{h \in \mathcal{H}} R(h) - R^* \right)}_{\text{approximation}}.$$

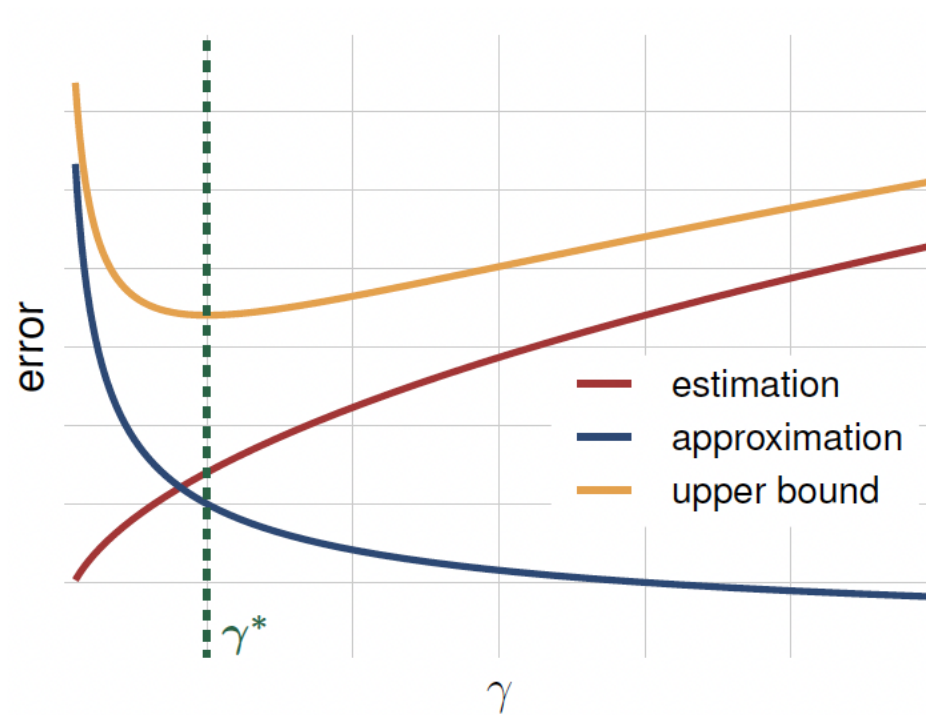


Estimation and approximation

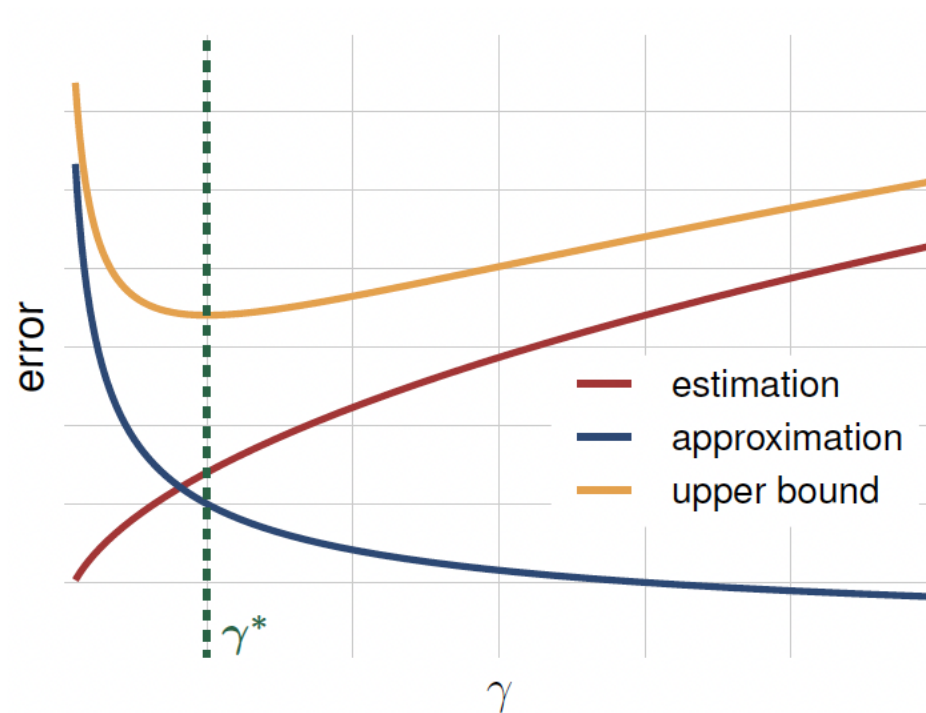
$$R(h) - R^* = \underbrace{\left(R(h) - \inf_{h \in \mathcal{H}} R(h) \right)}_{\text{estimation}} + \underbrace{\left(\inf_{h \in \mathcal{H}} R(h) - R^* \right)}_{\text{approximation}}.$$



Trade-off: estimation and approximation



Trade-off: estimation and approximation

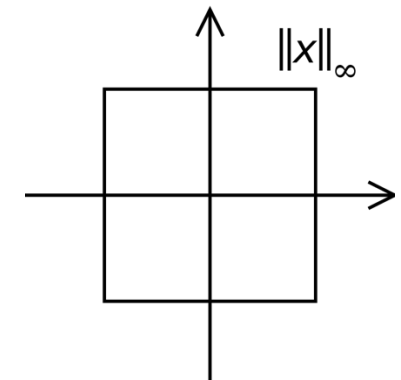
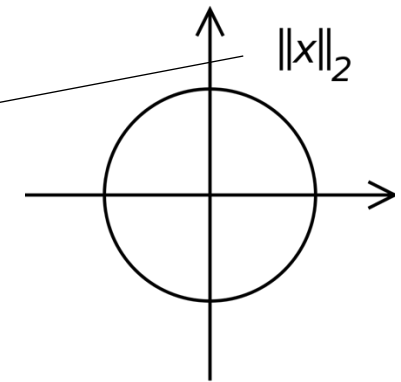
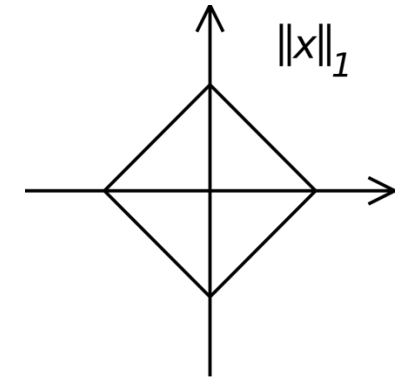


Q: how to control the richness of H?

Constrained problem

$$(CP) : \min_x f(x), \text{ s.t. } h(x) \leq b,$$

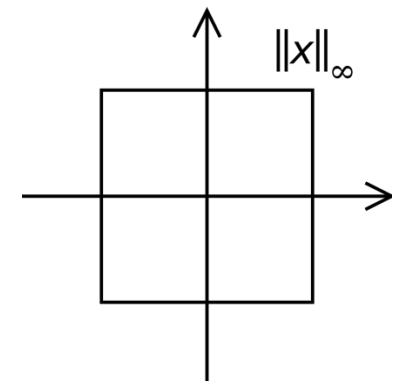
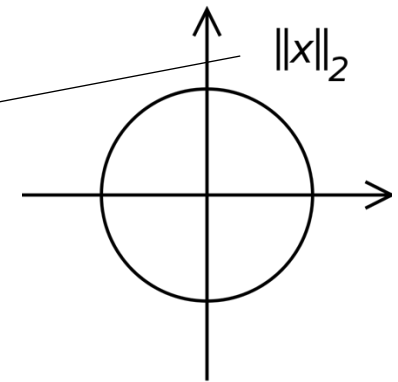
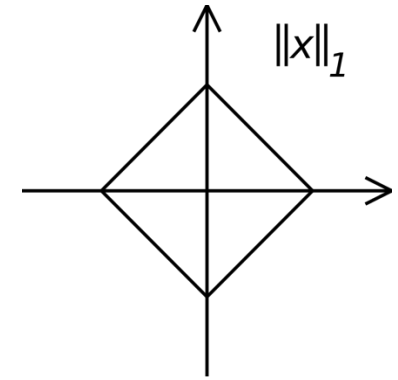
$$h(x) = \|x\|_2^2 \\ = x_1^2 + x_2^2 + \dots + x_n^2$$



Regularized problem

$$(RP) : \min_x f(x) + \lambda h(x),$$

$$h(x) = \|x\|_2^2 \\ = x_1^2 + x_2^2 + \dots + x_n^2$$



Regularized problem

$$(RP) : \min_x f(x) + \lambda h(x),$$

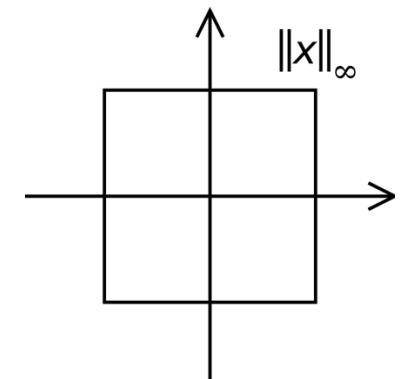
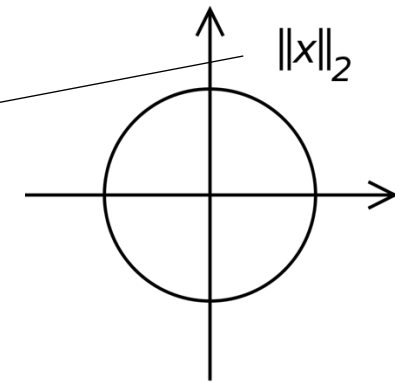
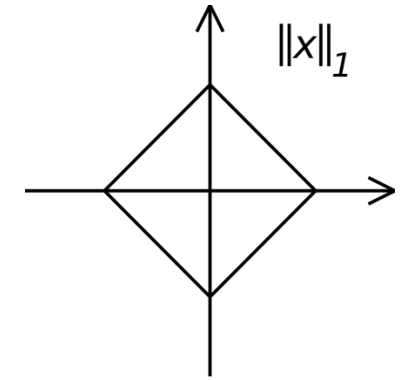
$$h(x) = \|x\|_2^2 \\ = x_1^2 + x_2^2 + \dots + x_n^2$$

Equivalence between (CP) and (RP):

$$\lambda \leftrightarrow b$$

We can find a b given λ such that:

Corresponding optimal solutions of (CP) and (RP) are identical



Empirical risk

Definition 2.2 (Empirical error) Given a hypothesis $h \in \mathcal{H}$, a target concept $c \in \mathcal{C}$, and a sample $S = (x_1, \dots, x_m)$ the empirical error or empirical risk of h is defined by

Training set

$$\hat{R}_S(h) = \frac{1}{m} \sum_{i=1}^m 1_{h(x_i) \neq c(x_i)}. \quad (2.2)$$

Interpret: average mistakes a hypothesis h makes on a sample

Empirical risk

Definition 2.2 (Empirical error) Given a hypothesis $h \in \mathcal{H}$, a target concept $c \in \mathcal{C}$, and a sample $S = (x_1, \dots, x_m)$ the empirical error or empirical risk of h is defined by

Training set

$$\hat{R}_S(h) = \frac{1}{m} \sum_{i=1}^m 1_{h(x_i) \neq c(x_i)}. \quad (2.2)$$

stochastic version

$$\hat{R}_S(f) := \frac{1}{n} \sum_{i=1}^n \mathbb{I}_{f(x_i) \neq y_i}.$$

Interpret: average mistakes a hypothesis h makes on a sample

$$R(h) = \mathbb{P}_{(x,y) \sim \mathcal{D}} [h(x) \neq y] = \mathbb{E}_{(x,y) \sim \mathcal{D}} [1_{h(x) \neq y}].$$

risk (in population): not accessible

Empirical risk minimization

$$h_S^{\text{ERM}} = \operatorname{argmin}_{h \in \mathcal{H}} \widehat{R}_S(h) \longrightarrow \text{Empirical risk}$$

Empirical risk minimization

$$R(h) - R^* = \underbrace{\left(R(h) - \inf_{h \in \mathcal{H}} R(h) \right)}_{\text{estimation}} + \underbrace{\left(\inf_{h \in \mathcal{H}} R(h) - R^* \right)}_{\text{approximation}}.$$

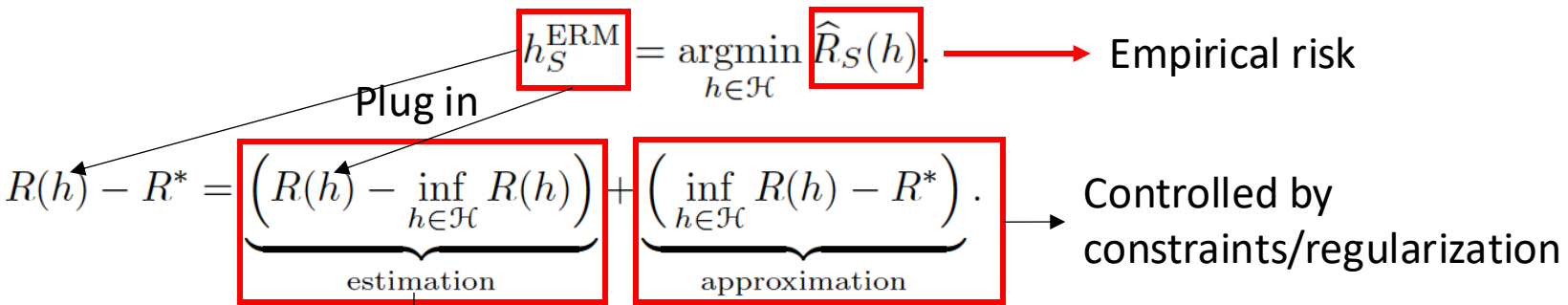
$h_S^{\text{ERM}} = \operatorname{argmin}_{h \in \mathcal{H}} \widehat{R}_S(h)$ → Empirical risk

Plug in

→ Controlled by constraints/regularization

Q: How to control?

Empirical risk minimization



Q: How to control?

Proposition 4.1 For any sample S , the following inequality holds for the hypothesis returned by ERM:

$$\mathbb{P} \left[R(h_S^{\text{ERM}}) - \inf_{h \in \mathcal{H}} R(h) > \epsilon \right] \leq \mathbb{P} \left[\sup_{h \in \mathcal{H}} |R(h) - \widehat{R}_S(h)| > \frac{\epsilon}{2} \right]. \quad (4.3)$$

Corollary 3.19 (VC-dimension generalization bounds) Let \mathcal{H} be a family of functions taking values in $\{-1, +1\}$ with VC-dimension d . Then, for any $\delta > 0$, with probability at least $1 - \delta$, the following holds for all $h \in \mathcal{H}$:

$$R(h) \leq \widehat{R}_S(h) + \sqrt{\frac{2d \log \frac{em}{d}}{m}} + \sqrt{\frac{\log \frac{1}{\delta}}{2m}} = O(\sqrt{1/m}) \quad (3.29)$$