# Lecture 4: Machine Learning II
## Principle of Learning
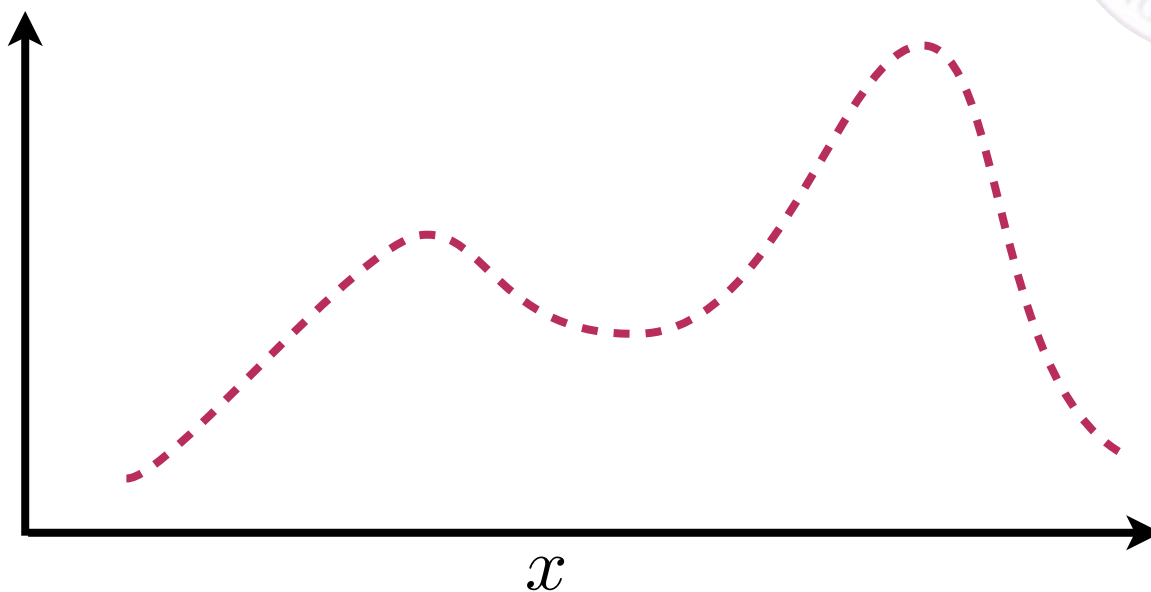
http://cs.nju.edu.cn/yuy/course_dm13ms.ashx

# The core of all the problems

infinite samples
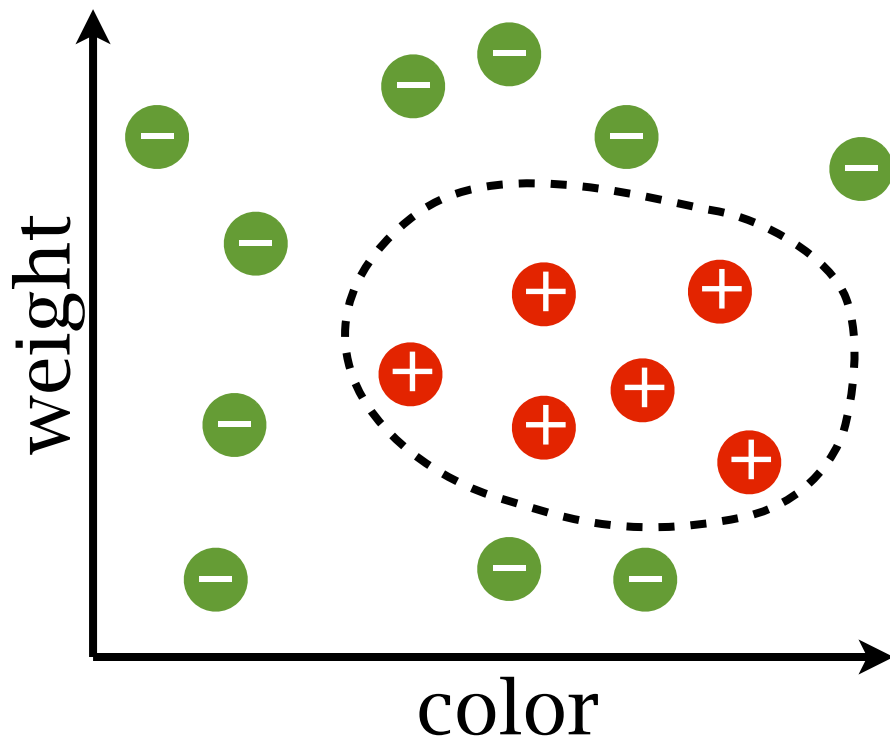
v.s.

finite samples

# Classification

**Features**: color, weight
**Label**: taste is sweet (positive/+) or not (negative/-)



(color, weight) $\rightarrow$ sweet ?

$$\mathcal{X} \rightarrow \{-1, +1\}$$

ground-truth function $f$

examples/training data:
$$\{(\boldsymbol{x}_1, y_1), \ldots, (\boldsymbol{x}_m, y_m)\}$$
$$y_i = f(\boldsymbol{x}_i)$$

learning: find an $f'$ that is <u>close</u> to $f$

# Classification

what can be observed:

on examples/training data:
$$\{(\boldsymbol{x}_1, y_1), \ldots, (\boldsymbol{x}_m, y_m)\} \quad y_i = f(\boldsymbol{x}_i)$$

e.g. training error
$$\epsilon_t = \frac{1}{m} \sum_{i=1}^{m} I(h(\boldsymbol{x}_i) \neq y_i)$$

what is expected:

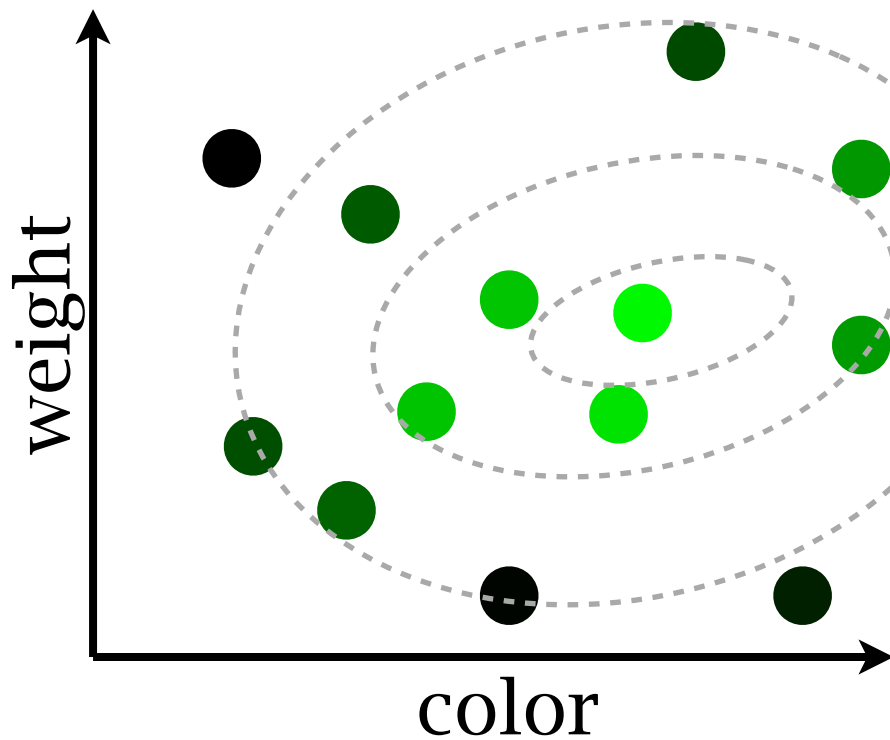over the whole distribution: generalization error
$$\epsilon_g = \mathbb{E}_x[I(h(\boldsymbol{x}) \neq f(\boldsymbol{x}))]$$
$$= \int_{\mathcal{X}} p(x) I(h(\boldsymbol{x}) \neq f(\boldsymbol{x}))] \mathrm{d}x$$

# Regression

**Features**: color, weight
**Label**: price [0,1]



(color, weight) $\rightarrow$ price
$$\mathcal{X} \rightarrow [0, +1]$$

ground-truth function $f$

examples/training data:
$$\{(\boldsymbol{x}_1, y_1), \ldots, (\boldsymbol{x}_m, y_m)\}$$
$$y_i = f(\boldsymbol{x}_i)$$

learning: <u>find</u> an $f'$ that is <u>close</u> to $f$

# Regression

what can be observed:

on examples/training data:
$$\{(\boldsymbol{x}_1, y_1), \ldots, (\boldsymbol{x}_m, y_m)\} \quad y_i = f(\boldsymbol{x}_i)$$

e.g. training mean square error/MSE

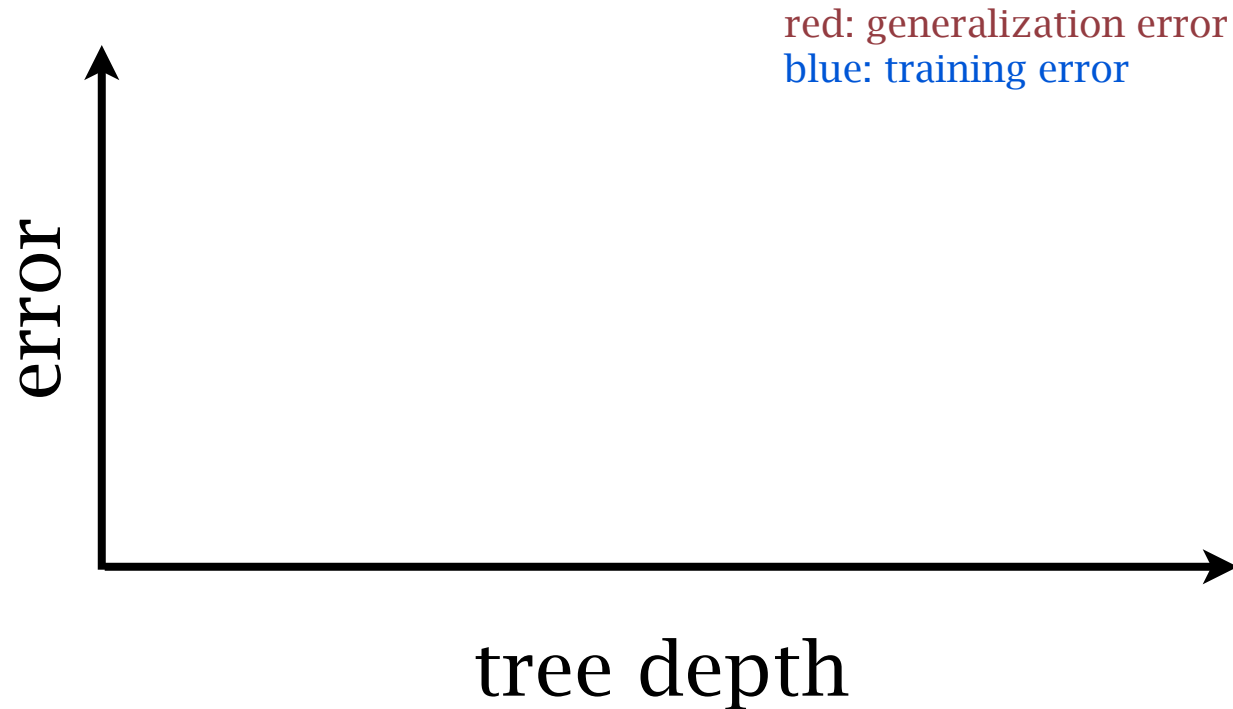$$\epsilon_t = \frac{1}{m} \sum_{i=1}^{m} (h(\boldsymbol{x}_i) - y_i)^2$$

what is expected:

over the whole distribution: generalization MSE

$$\epsilon_g = \mathbb{E}_x (h(\boldsymbol{x}) \neq f(\boldsymbol{x}))^2$$

$$= \int_{\mathcal{X}} p(x)(h(\boldsymbol{x}) - f(\boldsymbol{x}))^2 \mathrm{d}x$$

# The overfitting phenomena

-- the divergence between infinite and finite samples

red: generalization error
blue: training error



error

tree depth

# The overfitting phenomena

-- the divergence between infinite and finite samples

red: generalization error
blue: training error

error

tree depth

# The overfitting phenomena

-- the divergence between infinite and finite samples

red: generalization error
blue: training error

# The overfitting phenomena

-- the divergence between infinite and finite samples

red: generalization error
blue: training error



error

tree depth

# The overfitting phenomena

-- the divergence between infinite and finite samples

red: generalization error
blue: training error



error

tree depth

# The overfitting phenomena

-- the divergence between infinite and finite samples

red: generalization error
blue: training error



error

tree depth

# The overfitting phenomena

-- the divergence between infinite and finite samples

red: generalization error
blue: training error
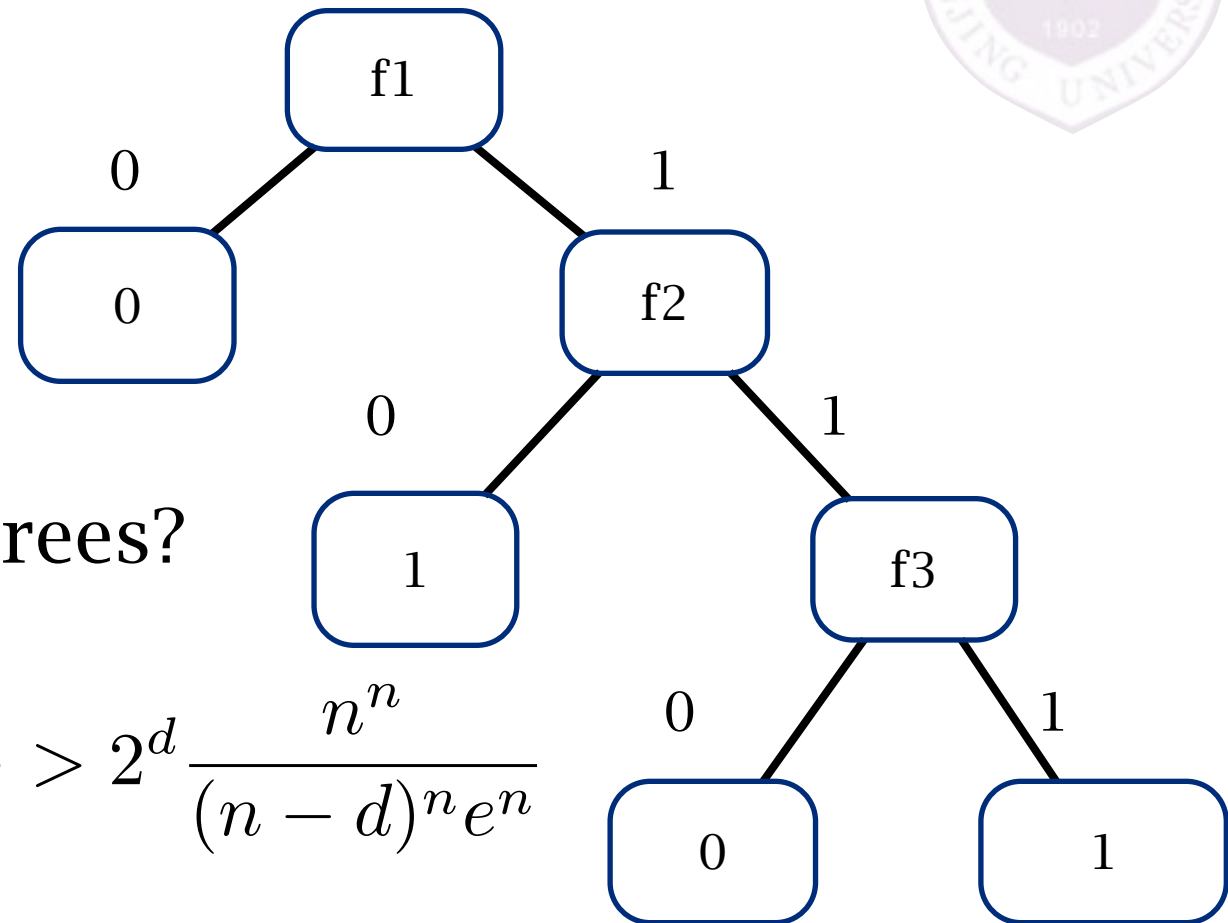


tree depth

*why tree depth?*

# Tree depth and the possibilities

features: $n$
feature type: binary
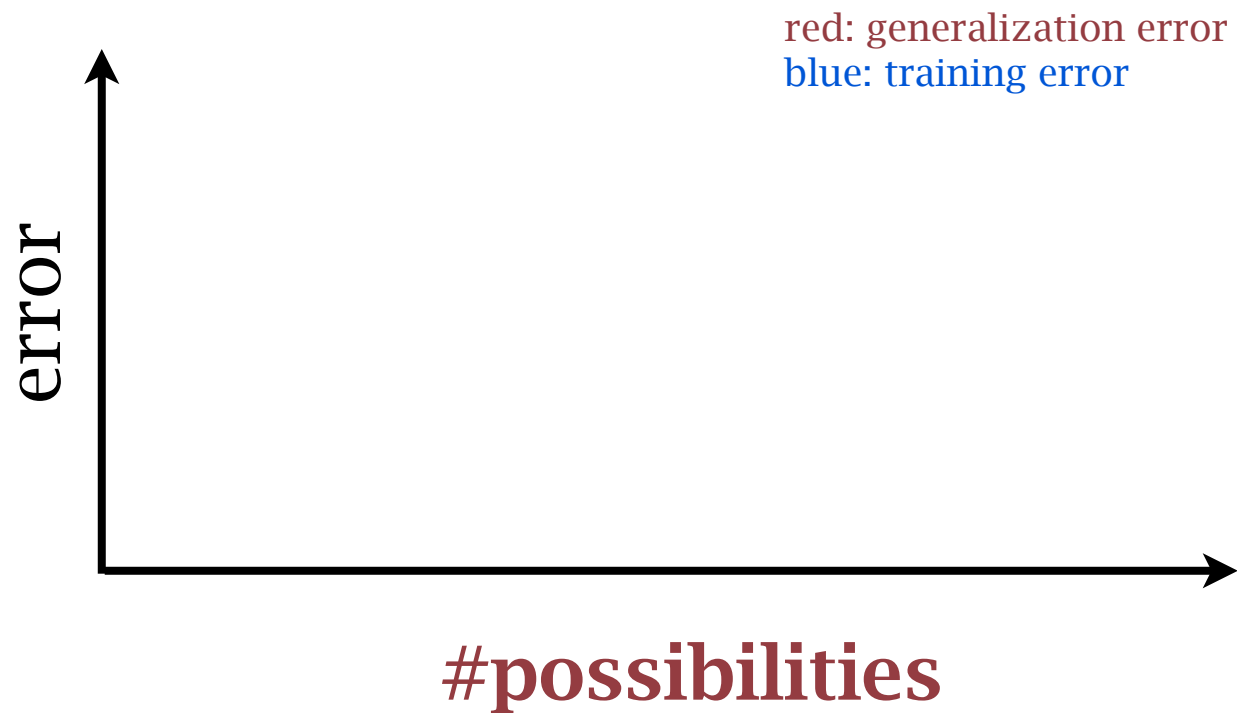depth: $d<n$

How many different trees?

one-branch: $2^d \dfrac{n!}{(n-d)!} > 2^d \dfrac{n^n}{(n-d)^n e^n}$

full-tree: $2^{2^d} \displaystyle\prod_{i=0}^{d-1} \dfrac{(n-i)!}{(n-d-i)!}$

the possibility of trees grows very fast with $d$

# The overfitting phenomena

-- the divergence between infinite and finite samples

red: generalization error
blue: training error



error

#possibilities

# The overfitting phenomena

-- the divergence between infinite and finite samples

red: generalization error
blue: training error

**#possibilities**

error

# The overfitting phenomena

-- the divergence between infinite and finite samples

red: generalization error
blue: training error



error

#possibilities

# The overfitting phenomena

-- the divergence between infinite and finite samples



red: generalization error
blue: training error

error

#possibilities

# The overfitting phenomena

-- the divergence between infinite and finite samples

red: generalization error
blue: training error



error

#possibilities

# The overfitting phenomena

-- the divergence between infinite and finite samples

red: generalization error
blue: training error



error

#possibilities

# The overfitting phenomena

-- the divergence between infinite and finite samples

red: generalization error
blue: training error



**#possibilities**

*why #possibilities?*

# The version space algorithm

an abstract view of learning algorithms



S: most specific hypothesis

G: most general hypothesis

version space: consistent hypotheses [Mitchell, 1997]

# The version space algorithm

an abstract view of learning algorithms



S: most specific hypothesis

G: most general hypothesis

version space: consistent hypotheses [Mitchell, 1997]

a conceptual algorithm:
1. for every example, remove the conflict boxes
2. find S in remaining boxes
3. find G in remaining boxes
4. output the mean of S and G
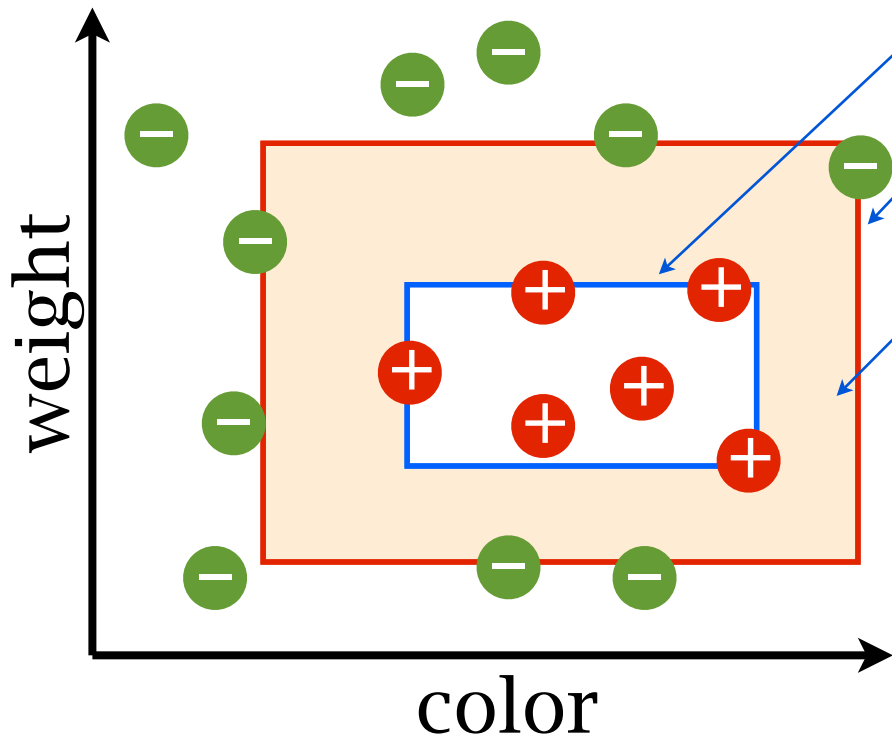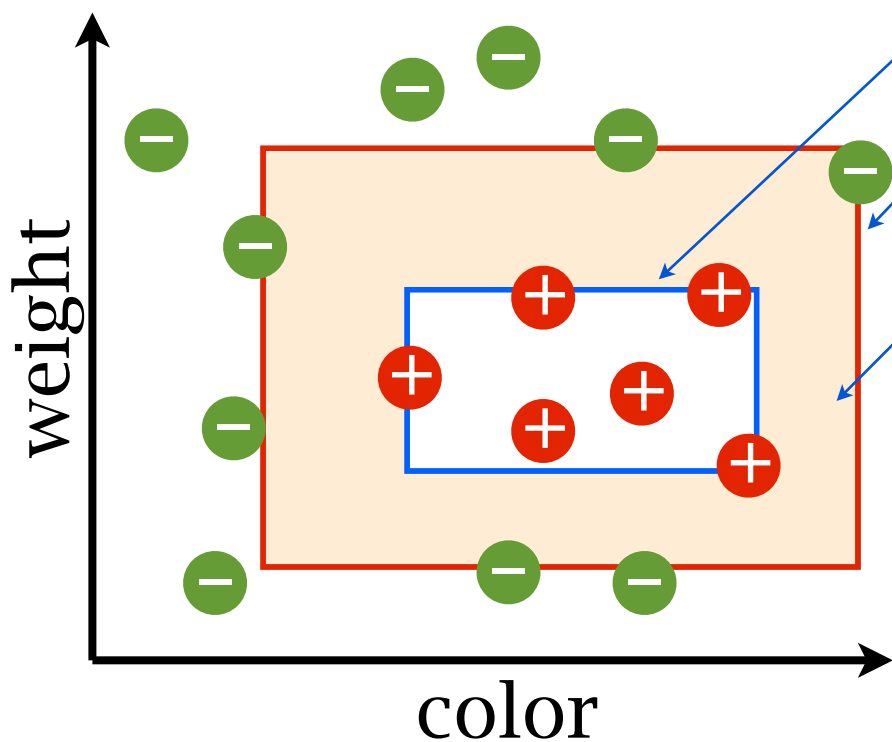
# The version space algorithm

an abstract view of learning algorithms



S: most specific hypothesis

G: most general hypothesis

version space: consistent hypotheses [Mitchell, 1997]

a conceptual algorithm:
1. for every example, remove the conflict boxes
2. find S in remaining boxes
3. find G in remaining boxes
4. output the mean of S and G

# The version space algorithm

an abstract view of learning algorithms



weight

color

S: most specific hypothesis

G: most general hypothesis

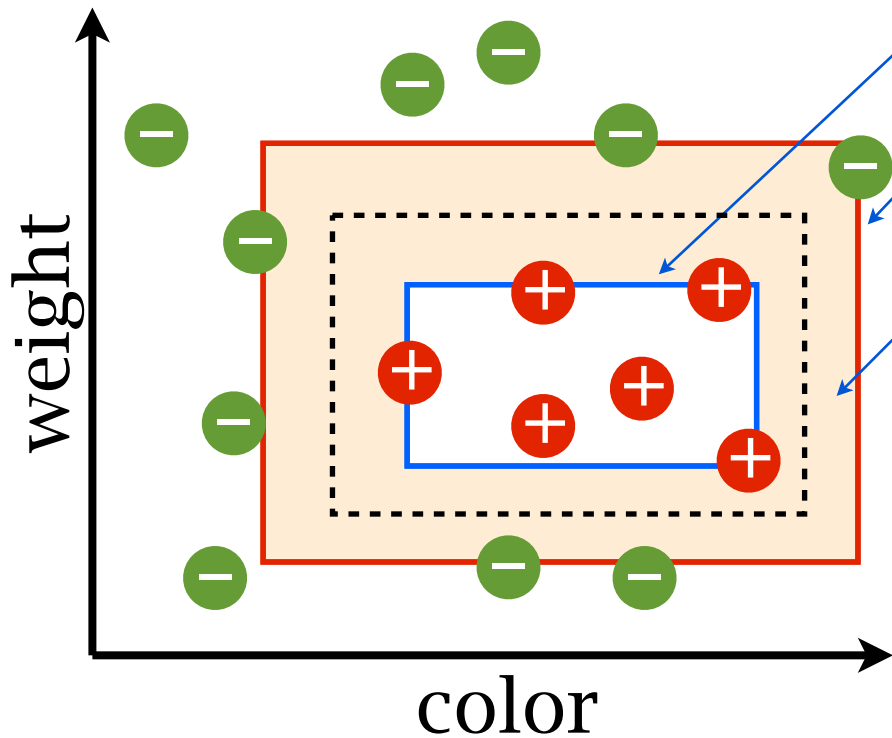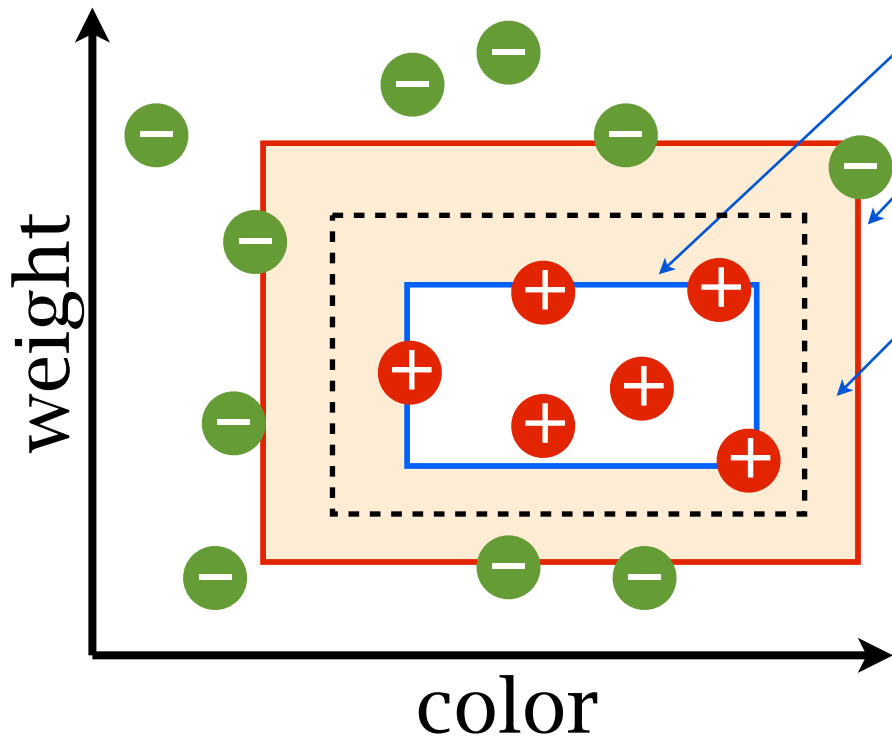version space: consistent hypotheses [Mitchell, 1997]

a conceptual algorithm:
1. for every example, remove the conflict boxes
2. find S in remaining boxes
3. find G in remaining boxes
4. output the mean of S and G

*selection a hypothesis according to learner's bias*

# The version space algorithm

an abstract view of learning algorithms

three components of a learning algorithm



#possibility ≈ hypothesis space size

# Theories

The i.i.d. assumption:

all training examples and future (test) examples are drawn *independently* from an *identical distribution*

unknown but fixed distribution $D$

bias-variance dilemma   (regression)

generalization bound    (classification)

# Bias-variance dilemma

Suppose we have 100 training examples
but there can be different 100 training examples

Start from the expected training MSE:

$$E_D[\epsilon_t] = E_D \left[ \frac{1}{m} \sum_{i=1}^{m} (h(\boldsymbol{x}_i) - y_i)^2 \right] = \frac{1}{m} \sum_{i=1}^{m} E_D \left[ (h(\boldsymbol{x}_i) - y_i)^2 \right]$$

(assume no noise)

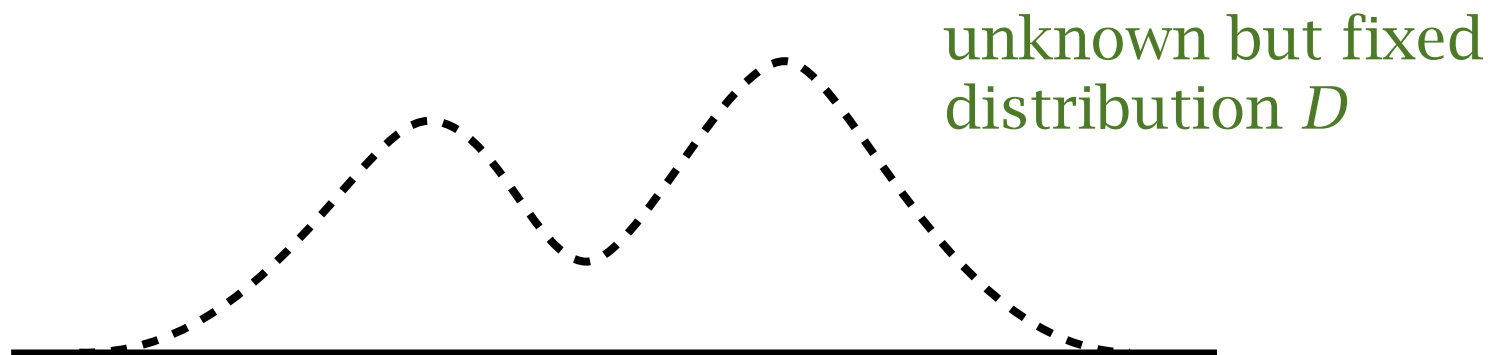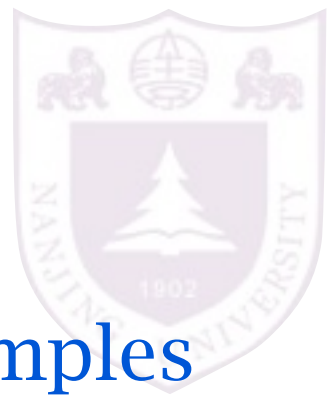$$= E_D \left[ (h(\boldsymbol{x}) - f(\boldsymbol{x}))^2 \right]$$

$$= E_D \left[ (h(\boldsymbol{x}) - E_D[h(\boldsymbol{x})] + E_D[h(\boldsymbol{x})] - f(\boldsymbol{x}))^2 \right]$$

$$= E_D \left[ (h(\boldsymbol{x}) - E_D[h(\boldsymbol{x})])^2 \right] + E_D \left[ (E_D[h(\boldsymbol{x})] - f(\boldsymbol{x}))^2 \right]$$

$$+ E_D \left[ 2(h(\boldsymbol{x}) - E_D[h(\boldsymbol{x})])(E_D[h(\boldsymbol{x})] - f(\boldsymbol{x})) \right]$$

$$= E_D \left[ (h(\boldsymbol{x}) - E_D[h(\boldsymbol{x})])^2 \right] + E_D \left[ (E_D[h(\boldsymbol{x})] - f(\boldsymbol{x}))^2 \right]$$

variance                                                bias^2

# Bias-variance dilemma

$$E_D\left[(h(\boldsymbol{x}) - E_D[h(\boldsymbol{x})])^2\right]$$
variance

$$E_D\left[(E_D[h(\boldsymbol{x})] - f(\boldsymbol{x}))^2\right]$$
bias^2



hypothesis space

# Bias-variance dilemma

$$E_D\left[(h(\boldsymbol{x}) - E_D[h(\boldsymbol{x})])^2\right]$$

variance

$$E_D\left[(E_D[h(\boldsymbol{x})] - f(\boldsymbol{x}))^2\right]$$

bias^2



hypothesis space

# Bias-variance dilemma

$$E_D\left[(h(\boldsymbol{x}) - E_D[h(\boldsymbol{x})])^2\right] \qquad E_D\left[(E_D[h(\boldsymbol{x})] - f(\boldsymbol{x}))^2\right]$$

variance                         bias^2

smaller hypothesis space

=>

smaller variance
but higher bias

$f$

hypothesis space

# Bias-variance dilemma

$$E_D\left[(h(\boldsymbol{x}) - E_D[h(\boldsymbol{x})])^2\right] \qquad E_D\left[(E_D[h(\boldsymbol{x})] - f(\boldsymbol{x}))^2\right]$$

$$\text{variance} \qquad\qquad \text{bias}\wedge 2$$

smaller hypothesis space
=>
smaller variance
but higher bias

$f$

hypothesis space

# Bias-variance dilemma

$$E_D\left[(h(\boldsymbol{x}) - E_D[h(\boldsymbol{x})])^2\right] \qquad E_D\left[(E_D[h(\boldsymbol{x})] - f(\boldsymbol{x}))^2\right]$$

$$\text{variance} \qquad\qquad\qquad \text{bias\^{}2}$$

high b
small v

balanced

low b
large v

red: generalization error
blue: training error

error

**|hypothesis space|**

# Bias-variance dilemma

$$E_D\left[(h(\boldsymbol{x}) - E_D[h(\boldsymbol{x})])^2\right] \qquad E_D\left[(E_D[h(\boldsymbol{x})] - f(\boldsymbol{x}))^2\right]$$

$$\text{variance} \qquad\qquad\qquad \text{bias\^2}$$

high b
small v

balanced

low b
large v

red: generalization error
blue: training error

error

|hypothesis space|

# Bias-variance dilemma

$$E_D\left[(h(\boldsymbol{x}) - E_D[h(\boldsymbol{x})])^2\right] \qquad E_D\left[(E_D[h(\boldsymbol{x})] - f(\boldsymbol{x}))^2\right]$$

variance                          bias^2

high b
small v          balanced          low b
                                   large v

error

red: generalization error
blue: training error

|hypothesis space|

# Bias-variance dilemma

$$E_D\left[(h(\boldsymbol{x}) - E_D[h(\boldsymbol{x})])^2\right] \qquad E_D\left[(E_D[h(\boldsymbol{x})] - f(\boldsymbol{x}))^2\right]$$

$$\text{variance} \qquad\qquad\qquad \text{bias\^{}2}$$



high b
small v

balanced

low b
large v

red: generalization error
blue: training error

error

|hypothesis space|

# Bias-variance dilemma

$$E_D\left[(h(\boldsymbol{x}) - E_D[h(\boldsymbol{x})])^2\right] \qquad E_D\left[(E_D[h(\boldsymbol{x})] - f(\boldsymbol{x}))^2\right]$$

$$\text{variance} \qquad\qquad\qquad\qquad \text{bias\textasciicircum{}2}$$

high b
small v

balanced

low b
large v

error

|hypothesis space|

red: generalization error
blue: training error

# Bias-variance dilemma

$$E_D\left[(h(\boldsymbol{x}) - E_D[h(\boldsymbol{x})])^2\right] \qquad E_D\left[(E_D[h(\boldsymbol{x})] - f(\boldsymbol{x}))^2\right]$$

$$\text{variance} \qquad\qquad\qquad \text{bias\textasciicircum 2}$$

high b
small v

balanced

low b
large v

red: generalization error
blue: training error

error

|hypothesis space|

# Overfitting and underfitting

training error v.s. hypothesis space size

# Overfitting and underfitting

training error v.s. hypothesis space size



linear functions: high training error, small space

$$\{y = a + bx \mid a, b \in \mathbb{R}\}$$
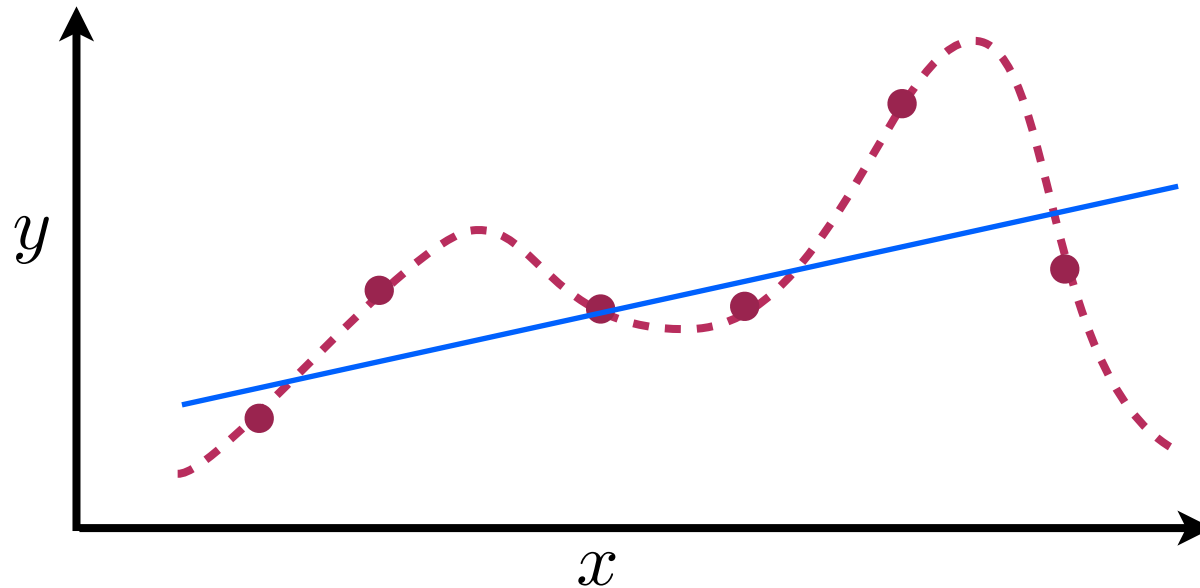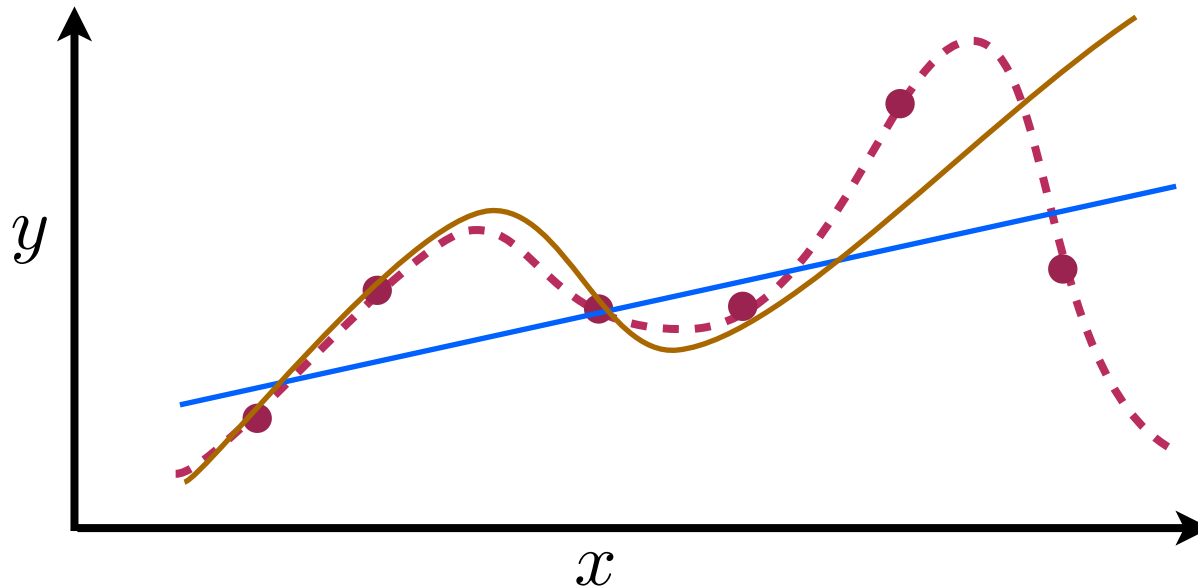
# Overfitting and underfitting

training error v.s. hypothesis space size



linear functions: high training error, small space
$$\{y = a + bx \mid a, b \in \mathbb{R}\}$$

higher polynomials: moderate training error, moderate space
$$\{y = a + bx + cx^2 + dx^3 \mid a, b, c, d \in \mathbb{R}\}$$

# Overfitting and underfitting

training error v.s. hypothesis space size



linear functions: high training error, small space
$$\{y = a + bx \mid a, b \in \mathbb{R}\}$$

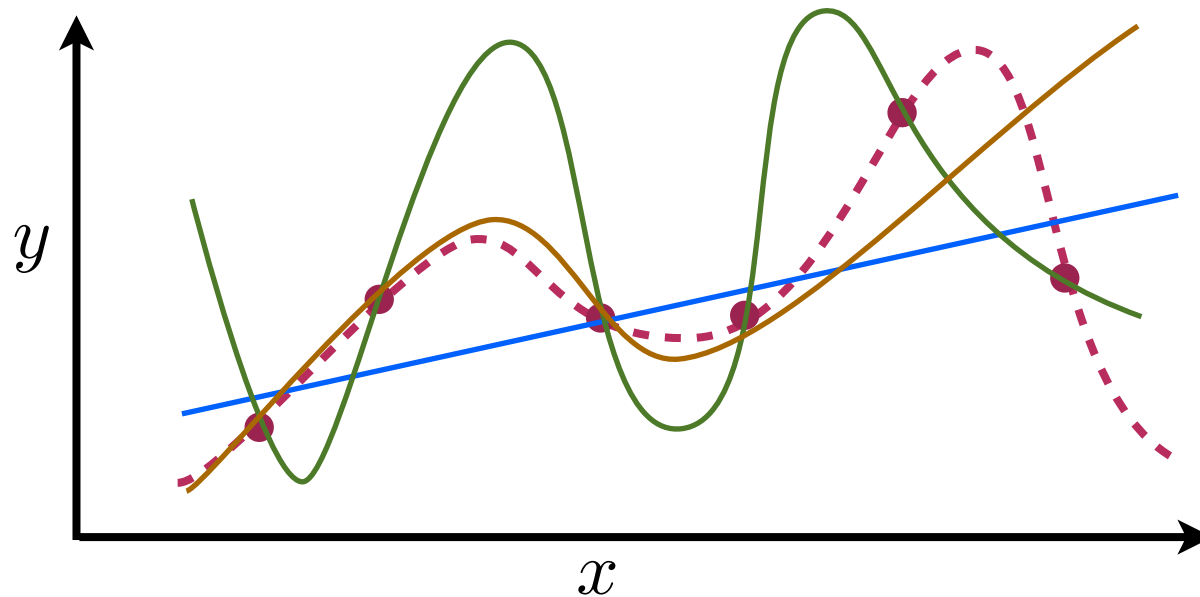higher polynomials: moderate training error, moderate space
$$\{y = a + bx + cx^2 + dx^3 \mid a, b, c, d \in \mathbb{R}\}$$

even higher order: no training error, large space
$$\{y = a + bx + cx^2 + dx^3 + ex^4 + fx^5 \mid a, b, c, d, e, f \in \mathbb{R}\}$$

# Generalization error

assume i.i.d. examples, and the ground-truth
hypothesis is a box

# Generalization error

assume i.i.d. examples, and the ground-truth hypothesis is a box



the error of picking a consistent hypothesis:

with probability at least $1 - \delta$

$$\epsilon_g < \frac{1}{m} \cdot (\ln |\mathcal{H}| + \ln \frac{1}{\delta})$$

# Generalization error

assume i.i.d. examples, and the ground-truth hypothesis is a box



the error of picking a consistent hypothesis:

with probability at least $1 - \delta$
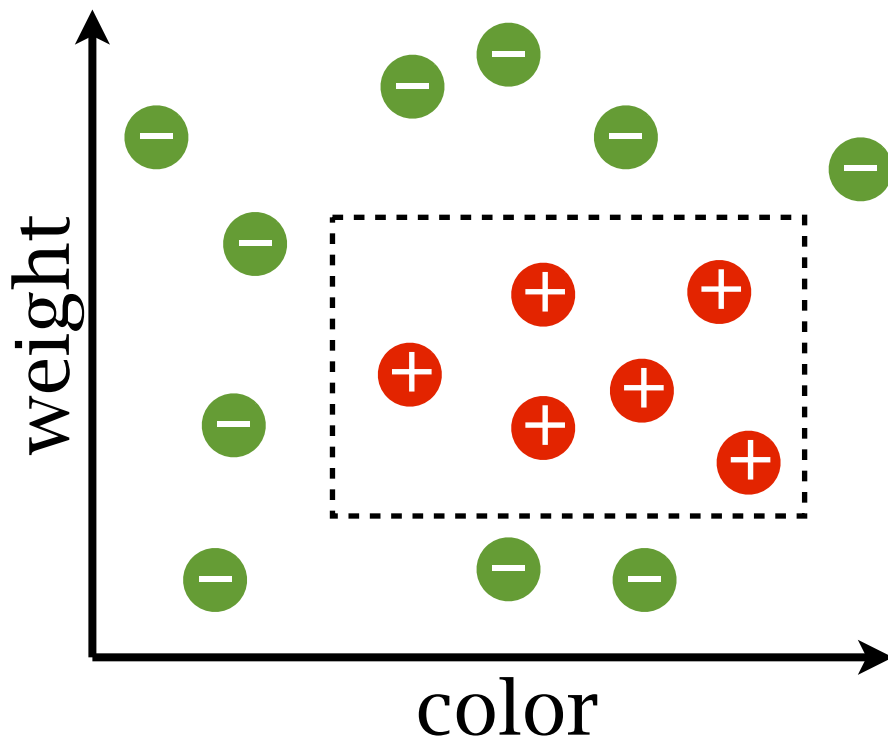
$$\epsilon_g < \frac{1}{m} \cdot \left( \ln |\mathcal{H}| + \ln \frac{1}{\delta} \right)$$

smaller generalization error:
▸ more examples
▸ smaller hypothesis space

# Generalization error

for one $h$

What is the probability of $\quad h$ is consistent

$\epsilon_g(h) \geq \epsilon$

assume $h$ is **bad**: $\epsilon_g(h) \geq \epsilon$

# Generalization error

for one $h$

What is the probability of
$$h \text{ is consistent}$$
$$\epsilon_g(h) \geq \epsilon$$

assume $h$ is **bad**: $\epsilon_g(h) \geq \epsilon$

$h$ is consistent with 1 example:

# Generalization error

for one $h$

What is the probability of $\quad$ $h$ is consistent

$\epsilon_g(h) \geq \epsilon$

assume $h$ is **bad**: $\epsilon_g(h) \geq \epsilon$

$h$ is consistent with 1 example:

$$P \leq 1 - \epsilon$$

# Generalization error

for one $h$

What is the probability of $\quad h$ is consistent

$\epsilon_g(h) \geq \epsilon$

assume $h$ is **bad**: $\epsilon_g(h) \geq \epsilon$

$h$ is consistent with 1 example:

$$P \leq 1 - \epsilon$$

$h$ is consistent with $m$ example:

# Generalization error

for one $h$

What is the probability of $h$ is consistent

$$\epsilon_g(h) \geq \epsilon$$

assume $h$ is **bad**: $\epsilon_g(h) \geq \epsilon$

$h$ is consistent with 1 example:

$$P \leq 1 - \epsilon$$

$h$ is consistent with **$m$** example:

$$P \leq (1 - \epsilon)^m$$

# Generalization error

$h$ is consistent with $m$ example:
$$P \leq (1 - \epsilon)^m$$

There are $k$ consistent hypotheses



...

# Generalization error

$h$ is consistent with $m$ example:

$$P \leq (1 - \epsilon)^m$$

There are $k$ consistent hypotheses

Probability of choosing a bad one:
$h_1$ is chosen and $h_1$ is bad $\quad P \leq (1 - \epsilon)^m$
$h_2$ is chosen and $h_2$ is bad $\quad P \leq (1 - \epsilon)^m$
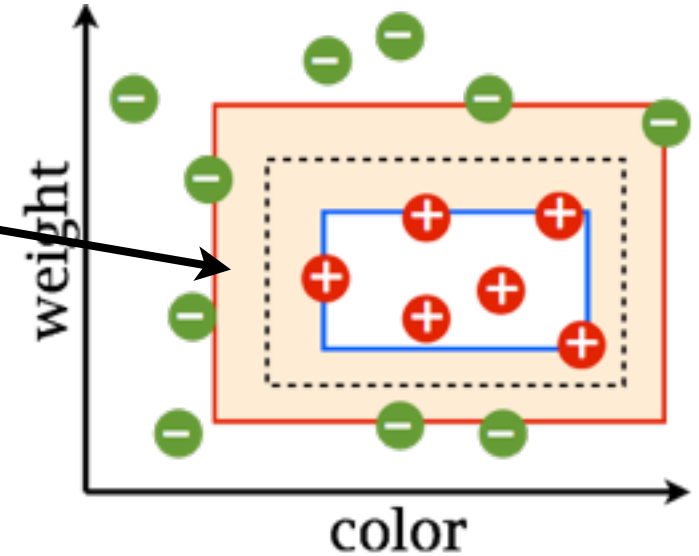...
$h_k$ is chosen and $h_k$ is bad $\quad P \leq (1 - \epsilon)^m$

# Generalization error

$h$ is consistent with $m$ example:

$$P \leq (1 - \epsilon)^m$$
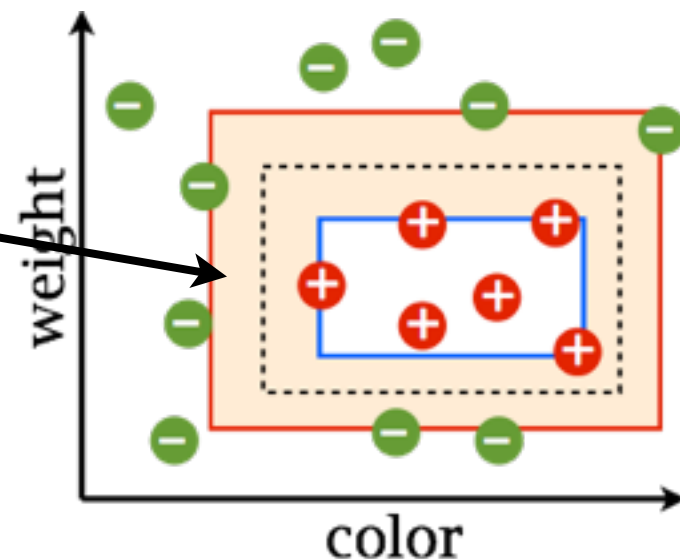
There are $k$ consistent hypotheses

Probability of choosing a bad one:
$h_1$ is chosen and $h_1$ is bad $\quad P \leq (1 - \epsilon)^m$
$h_2$ is chosen and $h_2$ is bad $\quad P \leq (1 - \epsilon)^m$
...
$h_k$ is chosen and $h_k$ is bad $\quad P \leq (1 - \epsilon)^m$

overall:
$\exists h$: $h$ can be chosen (consistent) but is bad

# Generalization error

$h_1$ is chosen and $h_1$ is bad $\quad P \leq (1 - \epsilon)^m$

$h_2$ is chosen and $h_2$ is bad $\quad P \leq (1 - \epsilon)^m$

...

$h_k$ is chosen and $h_k$ is bad $\quad P \leq (1 - \epsilon)^m$

overall:

$\exists h$: $h$ can be chosen (consistent) but is bad

# Generalization error

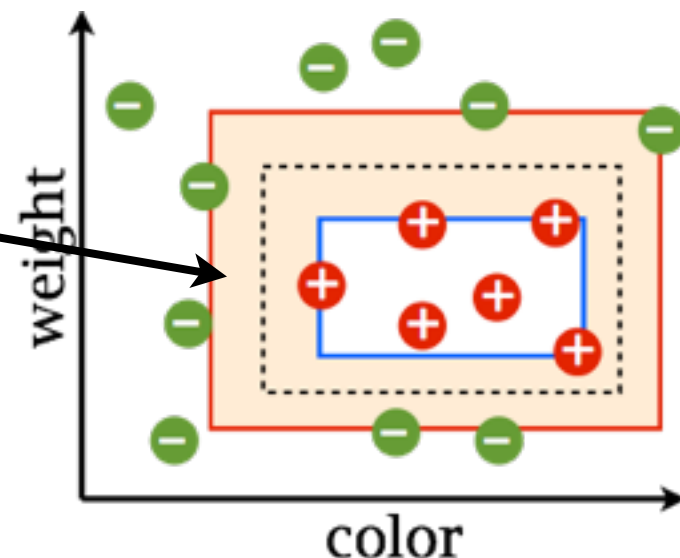$h_1$ is chosen and $h_1$ is bad $\quad P \le (1 - \epsilon)^m$

$h_2$ is chosen and $h_2$ is bad $\quad P \le (1 - \epsilon)^m$

...

$h_k$ is chosen and $h_k$ is bad $\quad P \le (1 - \epsilon)^m$

overall:

$\exists h$: $h$ can be chosen (consistent) but is bad

Union bound: $P(A \cup B) \le P(A) + P(B)$

# Generalization error

$h_1$ is chosen and $h_1$ is bad $\quad P \leq (1-\epsilon)^m$
$h_2$ is chosen and $h_2$ is bad $\quad P \leq (1-\epsilon)^m$
...
$h_k$ is chosen and $h_k$ is bad $\quad P \leq (1-\epsilon)^m$

overall:

$\exists h$: $h$ can be chosen (consistent) but is bad

Union bound: $P(A \cup B) \leq P(A) + P(B)$

$$P(\exists h \text{ is consistent but bad}) \leq k \cdot (1-\epsilon)^m \leq |\mathcal{H}| \cdot (1-\epsilon)^m$$

# Generalization error

$$P(\exists h \text{ is consistent but bad}) \leq k \cdot (1 - \epsilon)^m \leq |\mathcal{H}| \cdot (1 - \epsilon)^m$$

$$P(\epsilon_g \geq \epsilon) \leq |\mathcal{H}| \cdot (1 - \epsilon)^m$$

$$\epsilon_g < \frac{1}{m} \cdot (\ln |\mathcal{H}| + \ln \frac{1}{\delta})$$

# Generalization error

$$P(\exists h \text{ is consistent but bad}) \leq k \cdot (1 - \epsilon)^m \leq |\mathcal{H}| \cdot (1 - \epsilon)^m$$
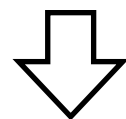
$$\Downarrow$$

$$P(\epsilon_g \geq \epsilon) \leq |\mathcal{H}| \cdot (1 - \epsilon)^m$$

$$\epsilon_g < \frac{1}{m} \cdot (\ln |\mathcal{H}| + \ln \frac{1}{\delta})$$

# Generalization error

$$P(\exists h \text{ is consistent but bad}) \leq k \cdot (1 - \epsilon)^m \leq |\mathcal{H}| \cdot (1 - \epsilon)^m$$
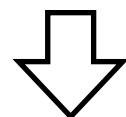
$$\Downarrow$$

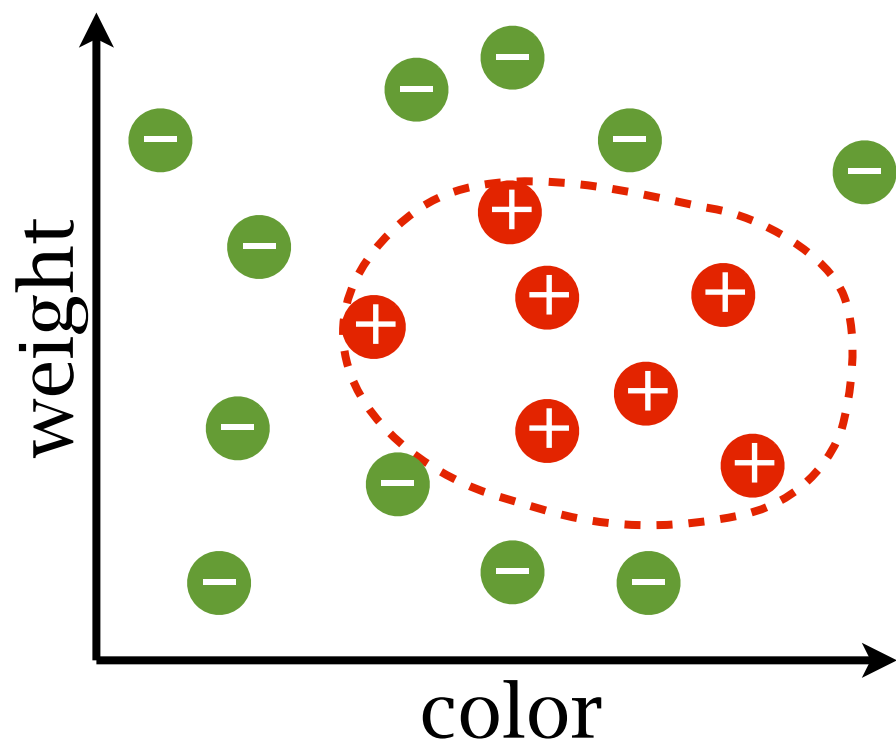$$P(\epsilon_g \geq \epsilon) \leq \frac{|\mathcal{H}| \cdot (1 - \epsilon)^m}{\delta}$$

with probability at least $1 - \delta$

$$\epsilon_g < \frac{1}{m} \cdot (\ln |\mathcal{H}| + \ln \frac{1}{\delta})$$

# Inconsistent hypothesis

What if the ground-truth hypothesis is NOT a box: non-zero training error
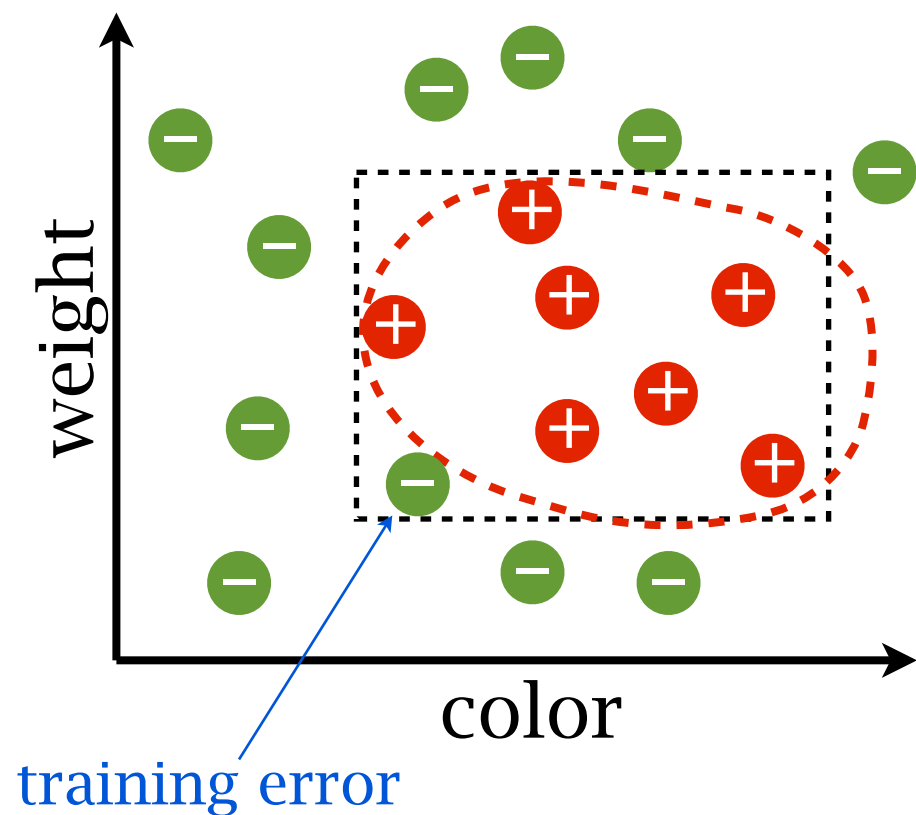
# Inconsistent hypothesis

What if the ground-truth hypothesis is NOT a box: non-zero training error

# Inconsistent hypothesis

What if the ground-truth hypothesis is NOT a box: non-zero training error



weight

color

training error

with probability at least $1 - \delta$

$$\epsilon_g < \epsilon_t + \sqrt{\frac{1}{m}\left(\ln |\mathcal{H}| + \ln \frac{1}{\delta}\right)}$$

# Inconsistent hypothesis

What if the ground-truth hypothesis is NOT a box: non-zero training error



training error

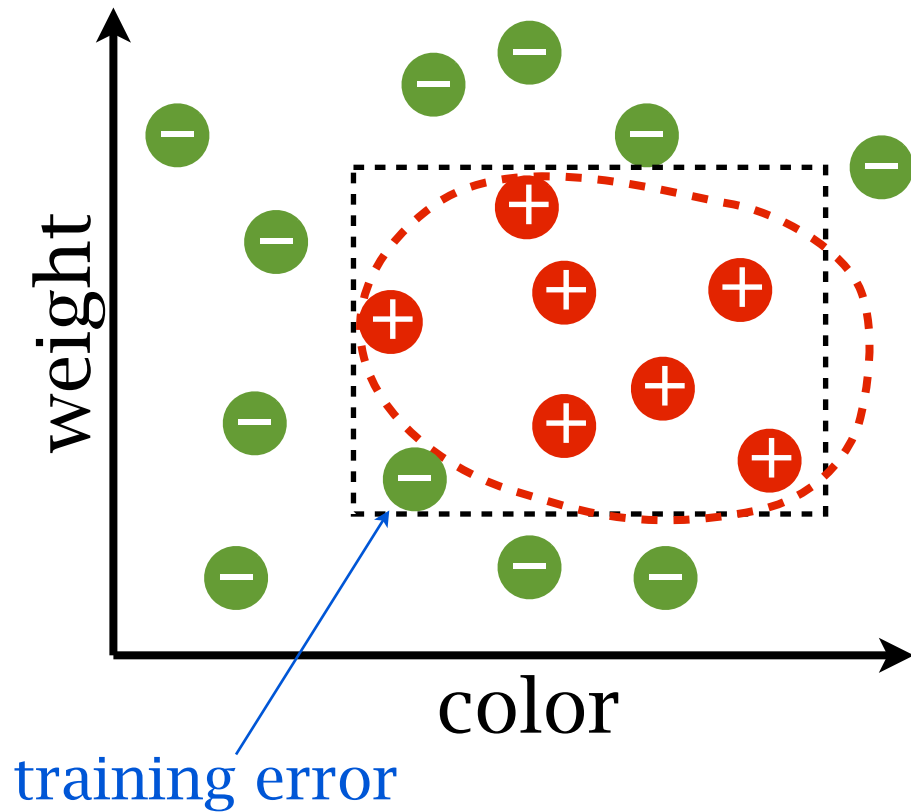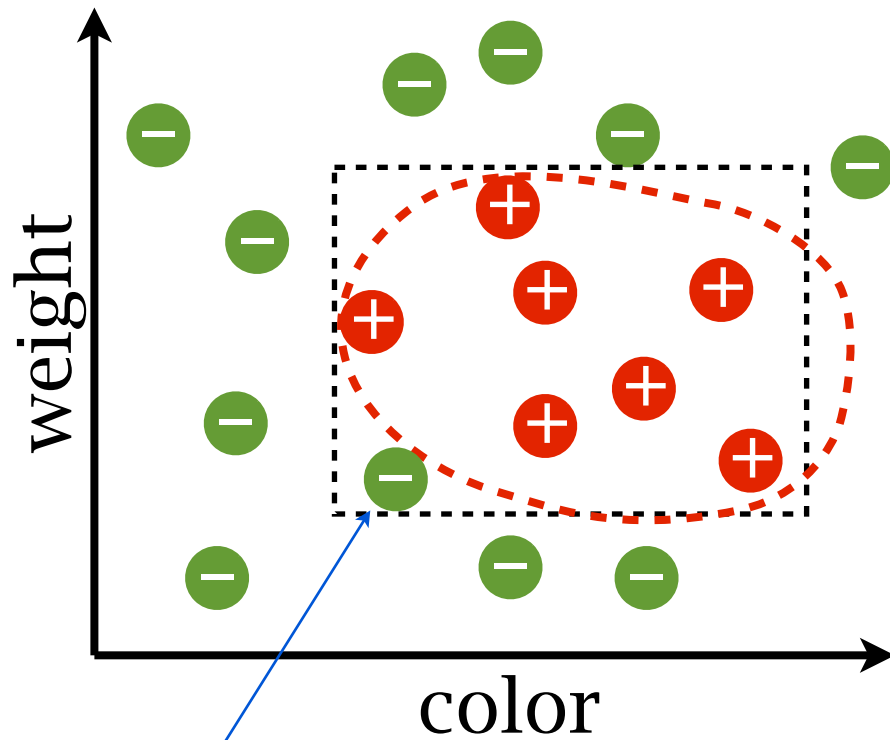with probability at least $1 - \delta$

$$\epsilon_g < \epsilon_t + \sqrt{\frac{1}{m}\left(\ln|\mathcal{H}| + \ln\frac{1}{\delta}\right)}$$

smaller generalization error:
- ▸ more examples
- ▸ smaller hypothesis space
- ▸ **smaller training error**

# Hoeffding's inequality

$X$ be an i.i.d. random variable
$X_1, X_2, \ldots, X_m$ be $m$ samples $\qquad X_i \in [b - a]$

$$\frac{1}{m} \sum_{i=1}^{m} X_i - \mathbb{E}[X] \leftarrow \text{ difference between sum and expectation}$$

$$P(\frac{1}{m} \sum_{i=1}^{m} X_i - \mathbb{E}[X] \geq \epsilon) \leq \exp\left(-\frac{2\epsilon^2 m}{(b-a)^2}\right)$$

# Generalization error

$$\text{for one } h$$

$$X_i = I(h(x_i) \neq f(x_i)) \in [0, 1]$$

$$\frac{1}{m} \sum_{i=1}^{m} X_i \rightarrow \epsilon_t(h) \qquad\qquad \mathbb{E}[X_i] \rightarrow \epsilon_g(h)$$

$$P(\epsilon_t(h) - \epsilon_g(h) \geq \epsilon) \leq \exp\left(-2\epsilon^2 m\right)$$

$$P(\epsilon_t - \epsilon_g \geq \epsilon)$$

$$\leq P(\exists h \in |\mathcal{H}| : \epsilon_t(h) - \epsilon_g(h) \geq \epsilon) \leq |\mathcal{H}| \exp\left(-2\epsilon^2 m\right)$$

$$\epsilon_g < \epsilon_t + \sqrt{\frac{1}{2m} \cdot \left(\ln |\mathcal{H}| + \ln \frac{1}{\delta}\right)}$$

# Generalization error

$$\text{for one } h$$

$$X_i = I(h(x_i) \neq f(x_i)) \in [0, 1]$$

$$\frac{1}{m} \sum_{i=1}^{m} X_i \to \epsilon_t(h) \qquad\qquad \mathbb{E}[X_i] \to \epsilon_g(h)$$

$$P(\epsilon_t(h) - \epsilon_g(h) \geq \epsilon) \leq \exp\left(-2\epsilon^2 m\right)$$

$$P(\epsilon_t - \epsilon_g \geq \epsilon)$$

$$\leq P(\exists h \in |\mathcal{H}| : \epsilon_t(h) - \epsilon_g(h) \geq \epsilon) \leq \frac{|\mathcal{H}| \exp\left(-2\epsilon^2 m\right)}{\delta}$$

$$\text{with probability at least } 1 - \delta$$

$$\epsilon_g < \epsilon_t + \sqrt{\frac{1}{2m} \cdot \left(\ln |\mathcal{H}| + \ln \frac{1}{\delta}\right)}$$

# Generalization error: Summary

assume i.i.d. examples
consistent hypothesis case:

with probability at least $1 - \delta$

$$\epsilon_g < \frac{1}{m} \cdot (\ln |\mathcal{H}| + \ln \frac{1}{\delta})$$

inconsistent hypothesis case:

with probability at least $1 - \delta$

$$\epsilon_g < \epsilon_t + \sqrt{\frac{1}{m}(\ln |\mathcal{H}| + \ln \frac{1}{\delta})}$$

generalization error:

number of examples $m$
training error $\epsilon_t$
hypothesis space complexity $\ln |\mathcal{H}|$

# PAC-learning

Probably approximately correct (PAC):

$$\epsilon_g < \epsilon_t + \sqrt{\frac{1}{2m} \cdot \left( \ln |\mathcal{H}| + \ln \frac{1}{\delta} \right)}$$

# PAC-learning

Probably approximately correct (PAC):

with probability at least $1 - \delta$

$$\epsilon_g < \epsilon_t + \sqrt{\frac{1}{2m} \cdot \left( \ln |\mathcal{H}| + \ln \frac{1}{\delta} \right)}$$

# PAC-learning

Probably approximately correct (PAC):

with probability at least $1 - \delta$

$$\epsilon_g < \epsilon_t + \sqrt{\frac{1}{2m} \cdot (\ln |\mathcal{H}| + \ln \frac{1}{\delta})}$$

**PAC-learnable:** [Valiant, 1984]

A concept class $\mathcal{C}$ is PAC-learnable if exists a learning algorithm $A$ such that for all $f \in \mathcal{C}$, $\epsilon > 0, \delta > 0$ and distribution $D$
$$P_D(\epsilon_g \leq \epsilon) \geq 1 - \delta$$
using $m = poly(1/\epsilon, 1/\delta)$ examples and polynomial time.

# PAC-learning

Probably approximately correct (PAC):

with probability at least $1 - \delta$

$$\epsilon_g < \epsilon_t + \sqrt{\frac{1}{2m} \cdot (\ln|\mathcal{H}| + \ln\frac{1}{\delta})}$$
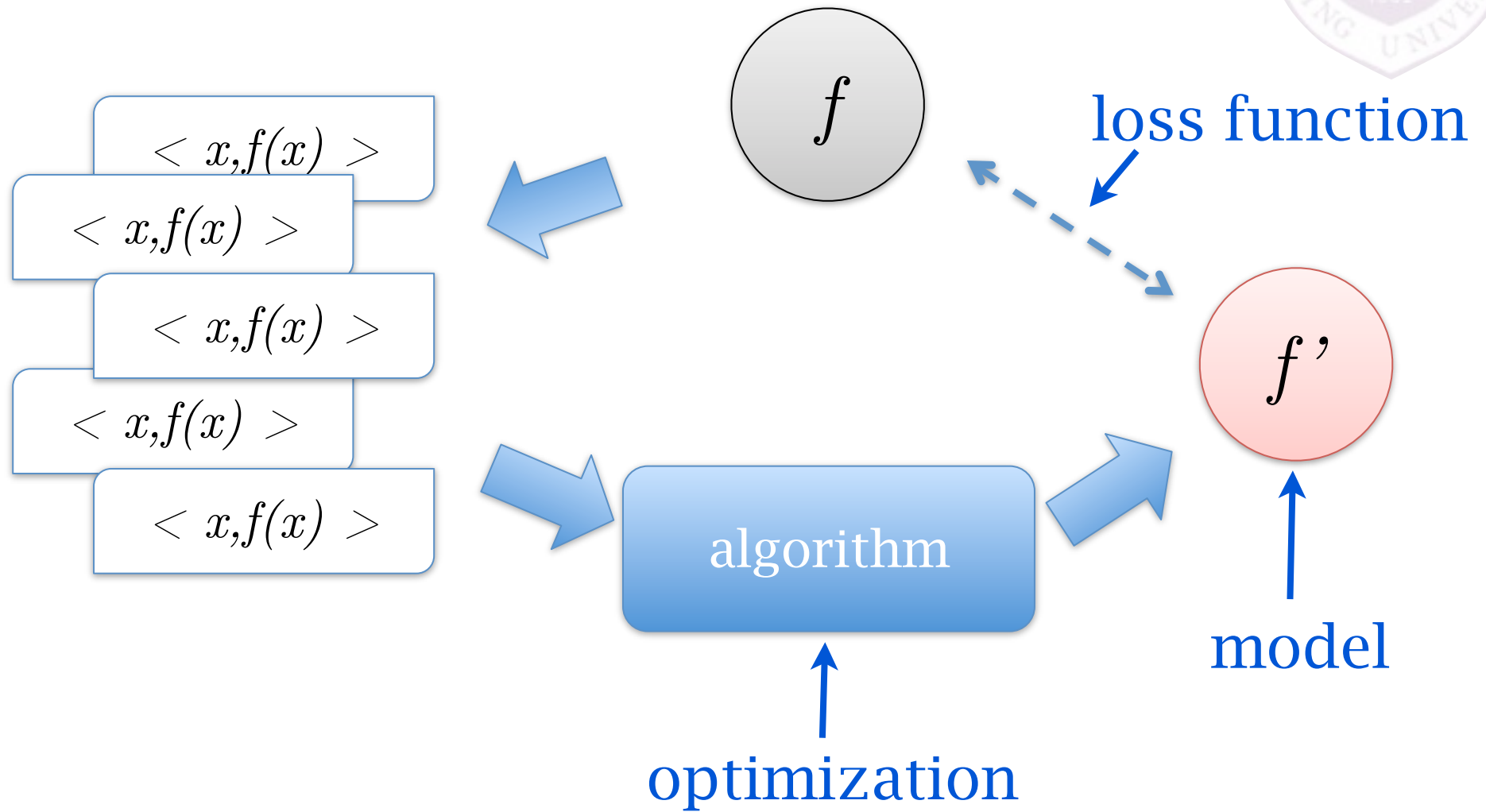
**PAC-learnable:** [Valiant, 1984]

A concept class $\mathcal{C}$ is PAC-learnable if exists a learning algorithm $A$ such that for all $f \in \mathcal{C}$, $\epsilon > 0, \delta > 0$ and distribution $D$

$$P_D(\epsilon_g \leq \epsilon) \geq 1 - \delta$$

using $m = poly(1/\epsilon, 1/\delta)$ examples and polynomial time.

Leslie Valiant
Turing Award (2010)
EATCS Award (2008)
Knuth Prize (1997)
Nevanlinna Prize (1986)

# Dimensions of modeling

监督学习的目标是否是最小化训练误差？

PAC-learning泛化界对于任意的潜在分布是否都成立？

以下两个多项式函数空间，哪一个的复杂度更高？
$$\mathcal{F}_1 = \{y = a + bx + cx^2 \mid a, b, c \in \mathbb{R}\}$$
$$\mathcal{F}_2 = \{y = a + ax + bx^2 + bx^3 + (a+b)x^4 \mid a, b \in \mathbb{R}\}$$

解释过配(overfitting)和欠配(underfitting)现象。

解释 Bias-Variance 困境