# Lecture 3:
## Supervised Learning

http://cs.nju.edu.cn/yuy/course_dm12.ashx
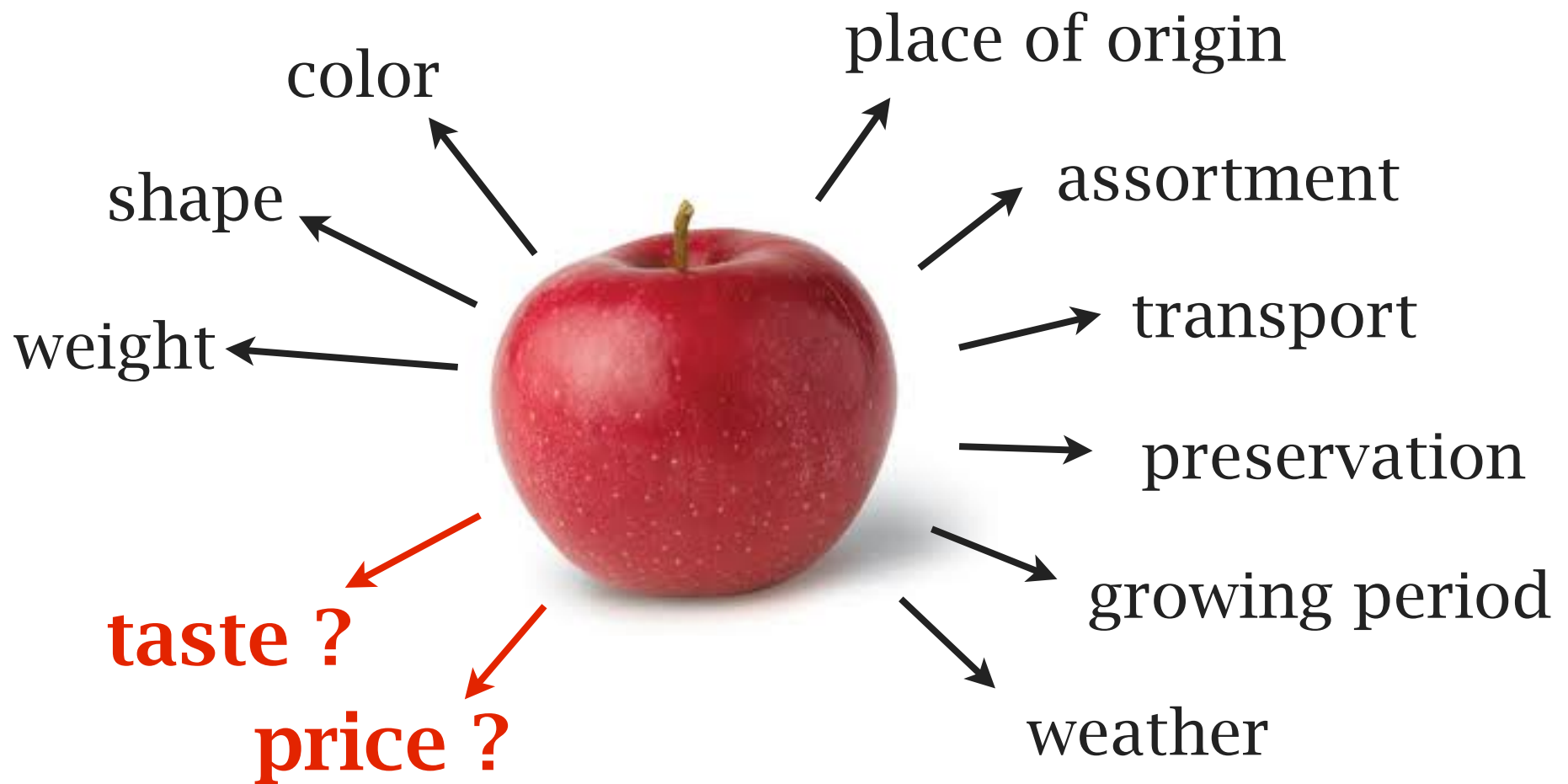
# Position

# The desire of prediction

# Predictive modeling

Find a relation between a set of variables (features) to target variables (labels).



color

shape

weight

taste ?

price ?

place of origin

assortment

transport

preservation

growing period

weather

# Supervised learning/inductive learning

Find a relation between a set of variables (features) to target variables (labels) *from finite examples*.

tasks

Classification: label is a nominal feature

Regression: label is a numerical feature

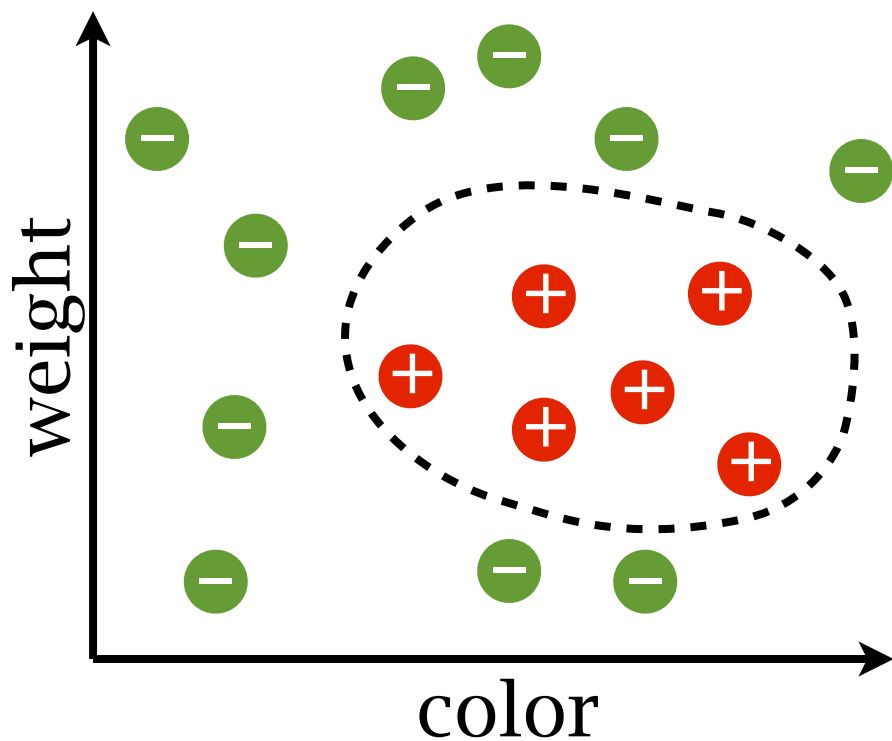Ranking: label is a ordinal feature

...

# Classification

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



(color, weight) → sweet ?

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

ground-truth function $f$

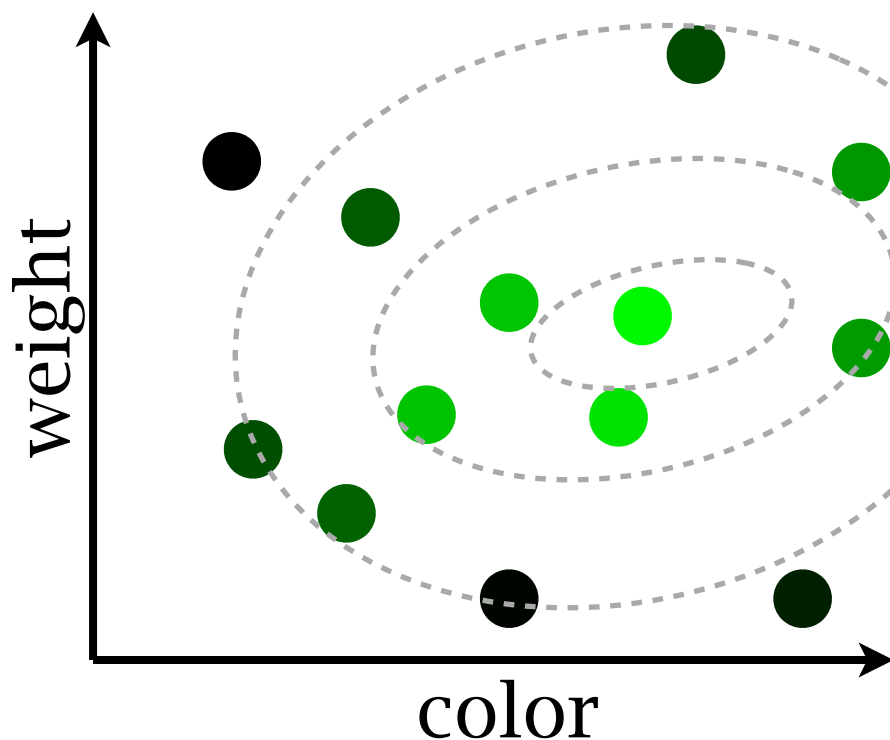examples/training data:
$$\{(\boldsymbol{x}_1, y_1), \ldots, (\boldsymbol{x}_m, y_m)\}$$
$$y_i = f(\boldsymbol{x}_i)$$

# Regression

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



(color, weight) → sweetness
$$\mathcal{X} \quad\quad \to [-1, +1]$$

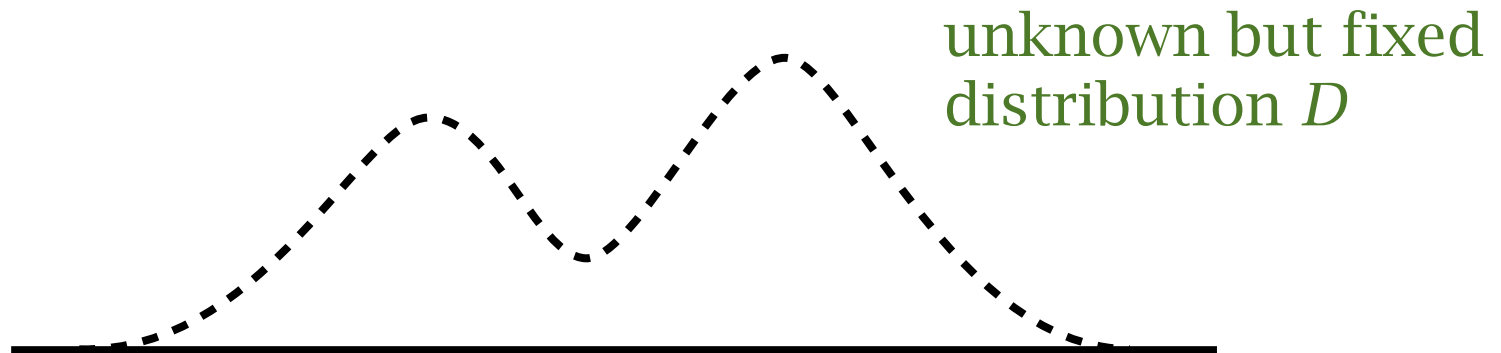ground-truth function $f$

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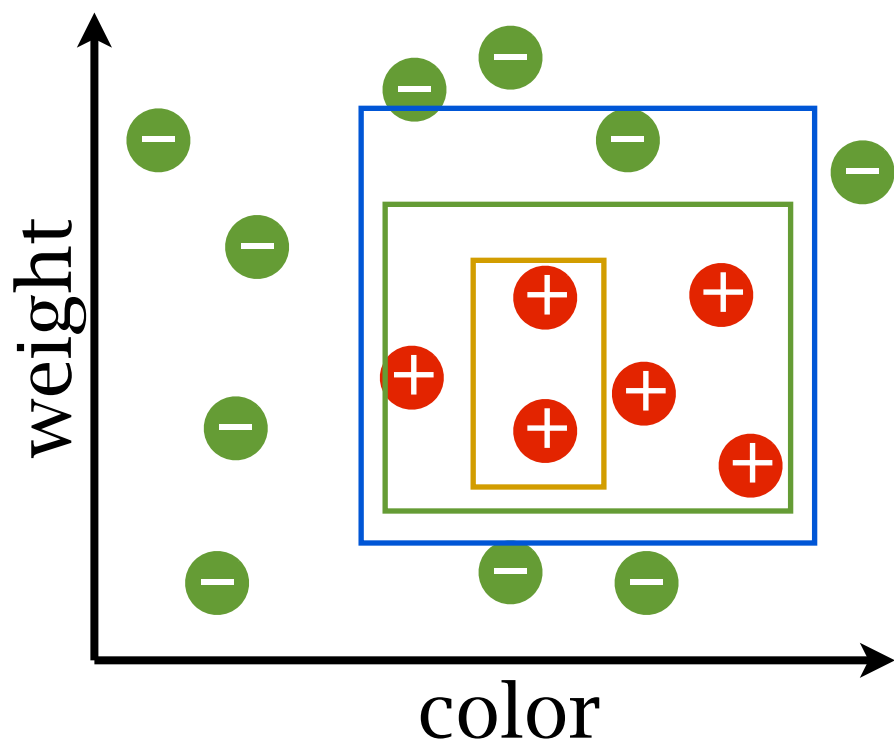$$y_i = f(\boldsymbol{x}_i)$$

weight

color

# I.I.D. assumption

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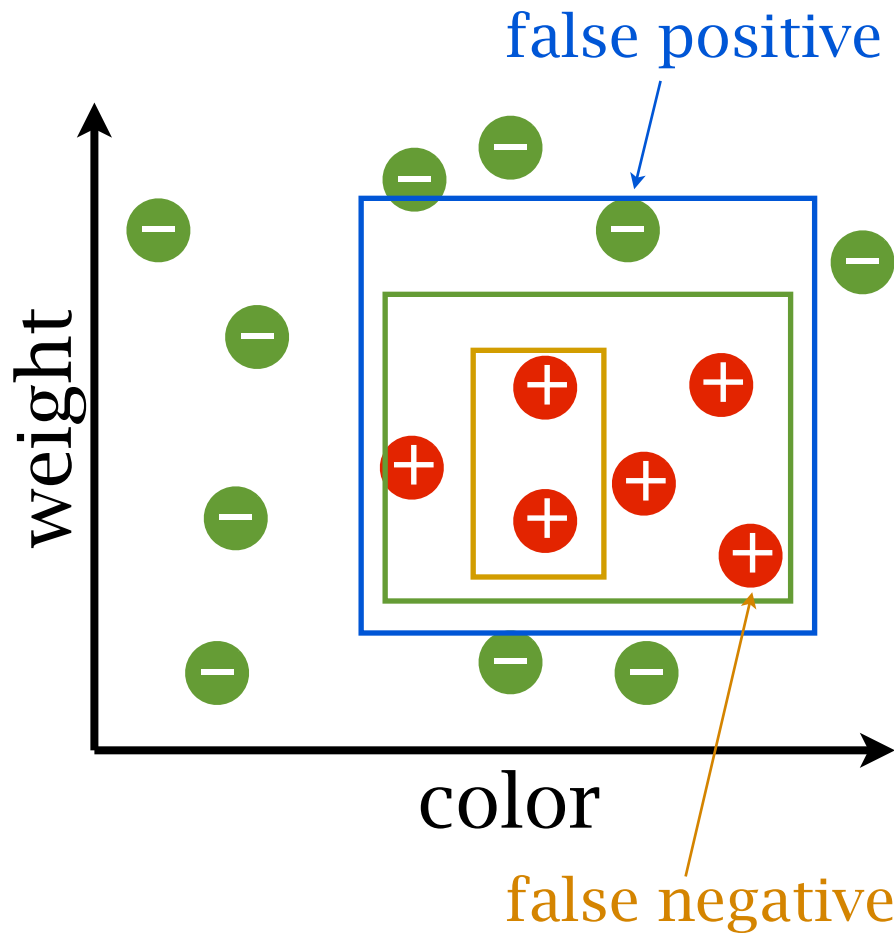an *identical distribution*



unknown but fixed
distribution *D*

box hypothesis class $\mathcal{H}$ contains all boxes

$h \in \mathcal{H}$ is a hypothesis

$$h(\boldsymbol{x}) = \begin{cases} +1, \text{if } x \text{ is inside the box} \\ -1, \text{if } x \text{ is outside the box} \end{cases}$$

# Training and generalization errors



false positive

weight

color

false negative

training error

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

generalization error

$$\epsilon_g = \mathbb{E}_x[I(h(\boldsymbol{x}) \neq f(\boldsymbol{x}))]$$

$$= \int_{\mathcal{X}} p(x) I(h(\boldsymbol{x}) \neq f(\boldsymbol{x}))]\mathrm{d}x$$

**find a hypothesis minimizes the generalization error**

# S, G, and the version space algorithm



S: most specific hypothesis

G: most general hypothesis
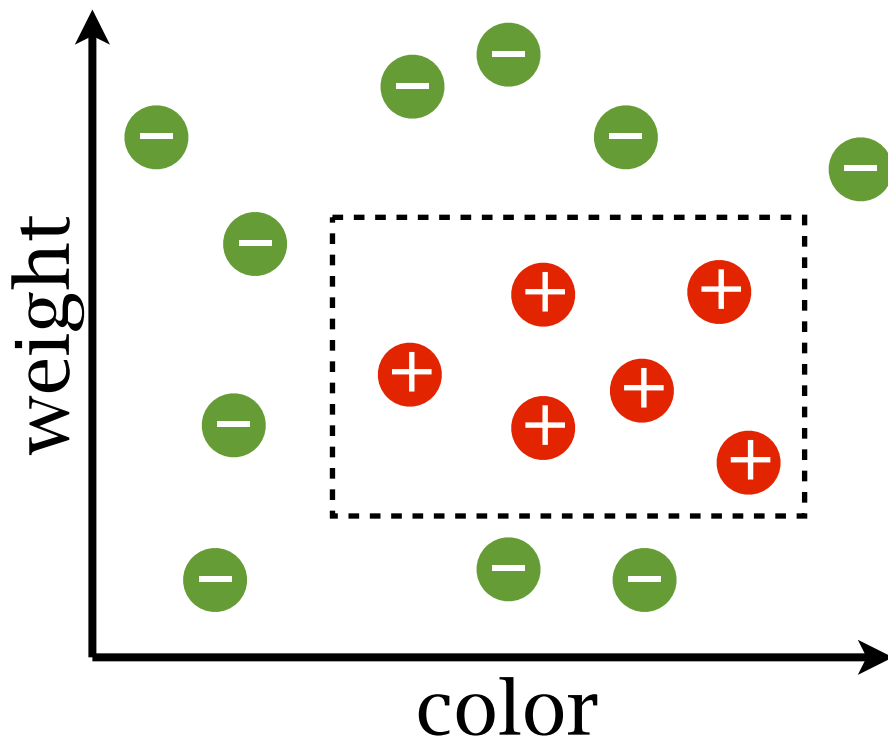
version space: consistent hypotheses [Mitchell, 1997]

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

# 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 $h$ 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

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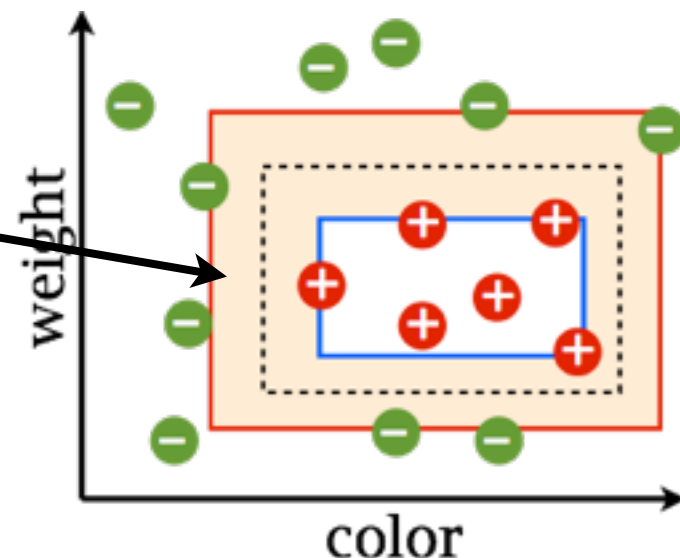$h_1$ is chosen and $h_1$ is bad $P \leq (1 - \epsilon)^m$
$h_2$ is chosen and $h_2$ is bad $P \leq (1 - \epsilon)^m$
...
$h_k$ is chosen and $h_k$ is bad $P \leq (1 - \epsilon)^m$

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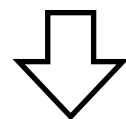
$\exists h$: $h$ can be chosen (consistent) but is bad

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

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

$$\Downarrow$$

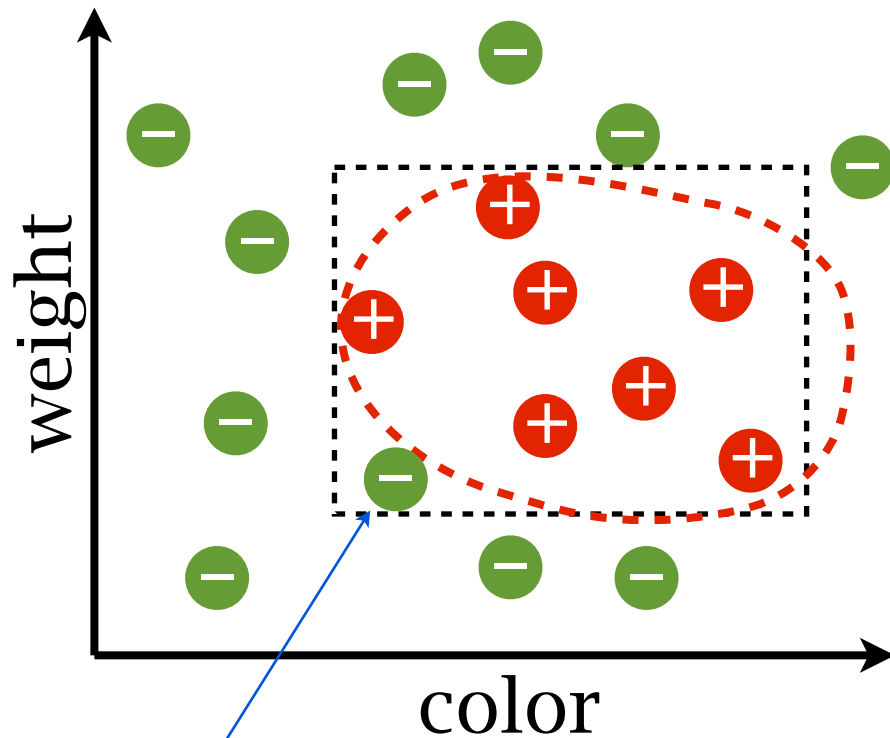$$P(\epsilon_g \geq \epsilon) \leq \frac{|\mathcal{H}| \cdot (1 - \epsilon)^m}{\delta}$$

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



training error

smaller generalization error:

with probability at least $1 - \delta$

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

▸ 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 $\qquad X_i \in [b - a]$

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

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

# Generalization error

$$\text{for one } h$$

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

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

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

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

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

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

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

# 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):

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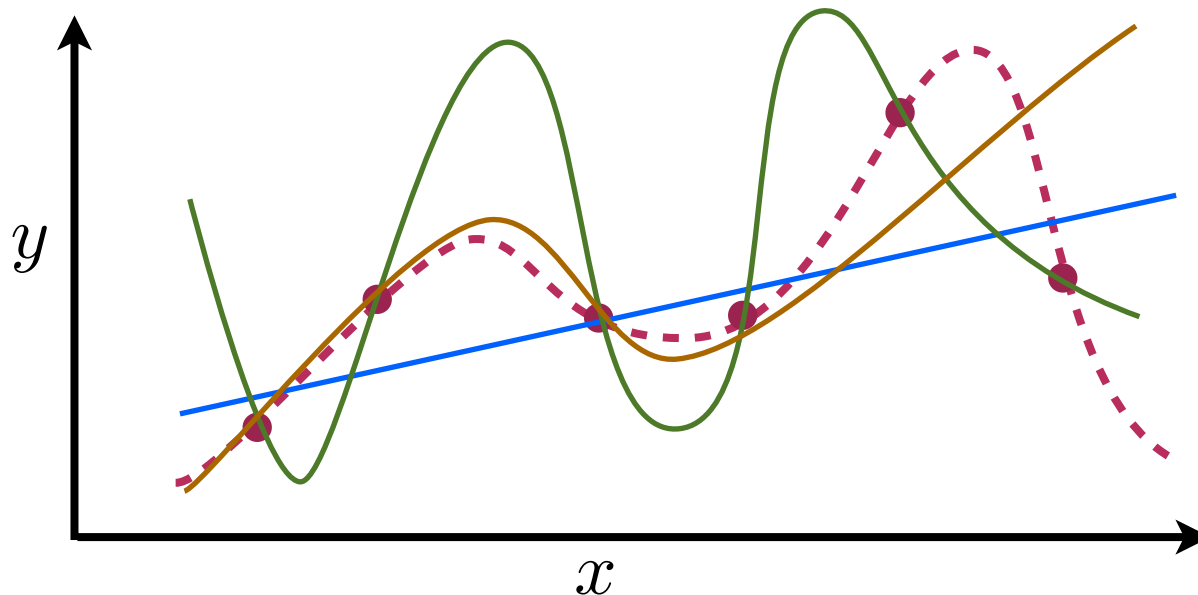
Leslie Valiant
Turing Award (2010)
EATCS Award (2008)
Knuth Prize (1997)
Nevanlinna Prize (1986)

# Overfitting and underfitting

training error v.s. hypothesis space size



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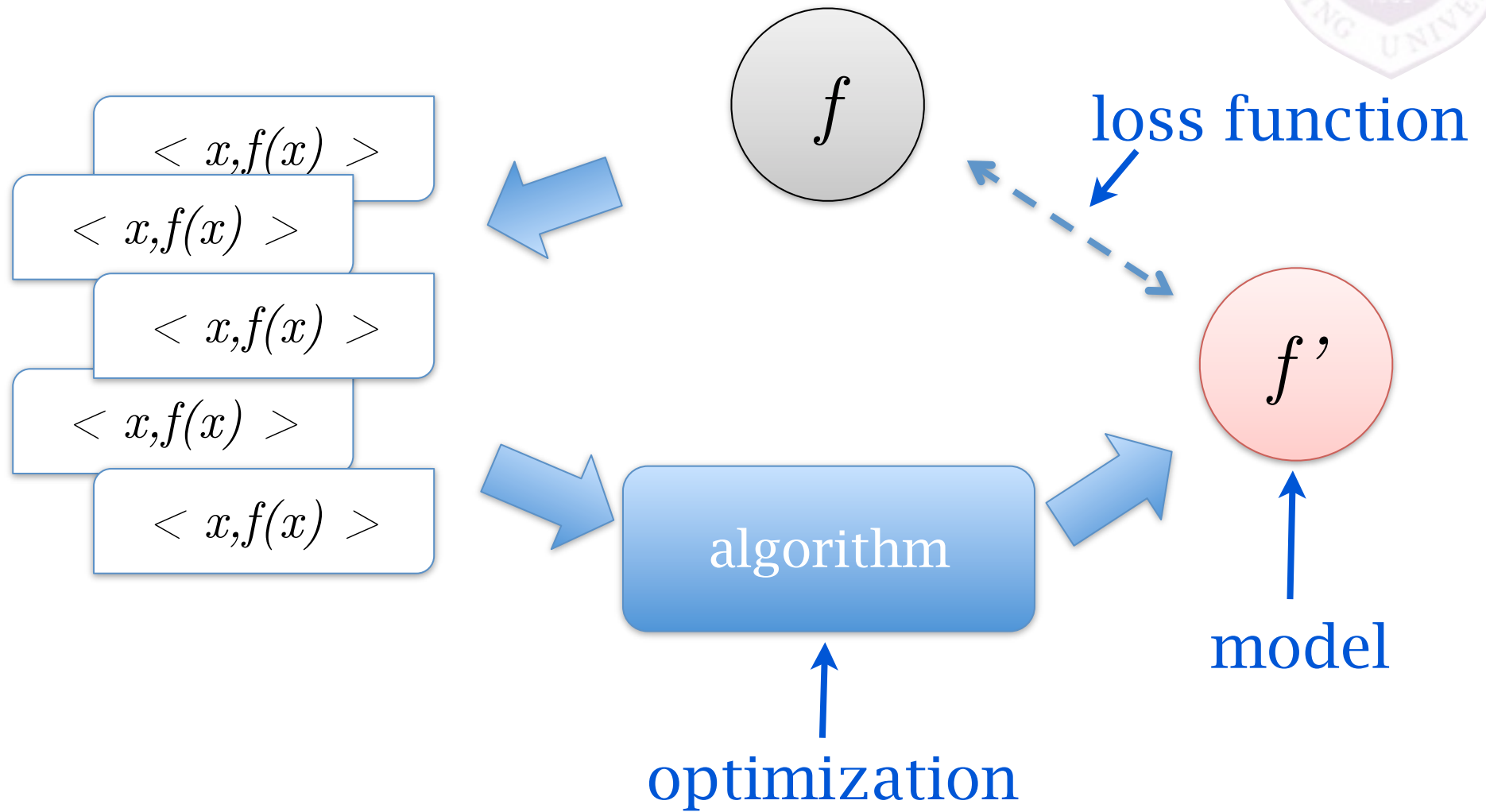$$\{y = a + bx \mid a, b \in \mathbb{R}\}$$

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
$$\{y = a + bx + cx^2 + dx^3 + ex^4 + fx^5 \mid a, b, c, d, e, f \in \mathbb{R}\}$$

# 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。