# Statistical Machine Learning

**The learning problem**

1 August 2018

# Outline

**1  Verification vs learning**

**2  Feasibility of learning**

**3  Error measures**

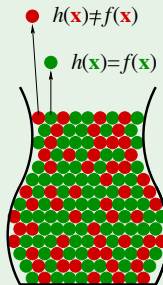**4  Noisy targets**

# Verification vs learning

Are we done?

Not so fast! $h$ is fixed.

For <u>this</u> $h$, $\nu$ generalizes to $\mu$.

'verification' of $h$, not **learning**

No guarantee $\nu$ will be small.

We need to **choose** from multiple $h$'s.



● $h(\mathbf{x}) \neq f(\mathbf{x})$

● $h(\mathbf{x}) = f(\mathbf{x})$
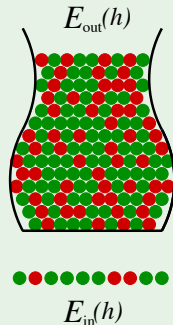
# Notation for learning

Both $\mu$ and $\nu$ depend on which hypothesis $h$

$\nu$ is 'in sample' denoted by $E_{\text{in}}(h)$

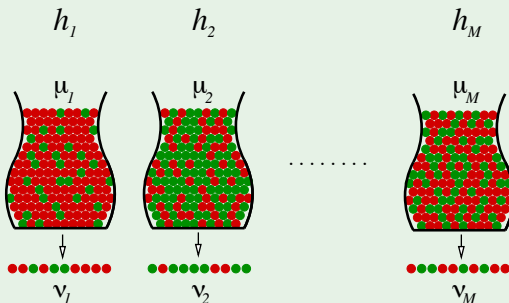$\mu$ is 'out of sample' denoted by $E_{\text{out}}(h)$

The Hoeffding inequality becomes:

$$\mathbb{P}\left[\,|E_{\text{in}}(h) - E_{\text{out}}(h)| > \epsilon\,\right] \;\leq\; 2e^{-2\epsilon^2 N}$$

$E_{\text{out}}(h)$



$E_{\text{in}}(h)$

# Multiple bins



## Multiple bins

Generalizing the bin model to more than one hypothesis:

$h_1$       $h_2$          $h_M$

$\mu_1$       $\mu_2$          $\mu_M$

. . . . . . . .

$\nu_1$       $\nu_2$          $\nu_M$

© Creator: Yaser Abu-Mostafa – LFD Lecture 2      10/17

# Notation with multiple bins



Notation with multiple bins

$h_1$ $\quad$ $h_2$ $\quad\quad\quad$ $h_M$

$E_{\text{out}}(h_1)$ $\quad$ $E_{\text{out}}(h_2)$ $\quad\quad\quad$ $E_{\text{out}}(h_M)$

. . . . . . . .

$E_{\text{in}}(h_1)$ $\quad$ $E_{\text{in}}(h_2)$ $\quad\quad\quad$ $E_{\text{in}}(h_M)$

# Coin analogy

## Hoefdding does not apply to multiple bins!

### Coin analogy

**Question:** If you toss a fair coin 10 times, what is the probability that you will get 10 heads?

**Answer:** $\approx 0.1\%$

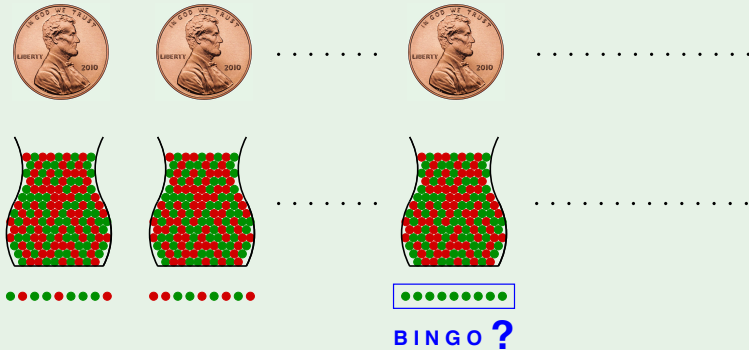**Question:** If you toss 1000 fair coins 10 times each, what is the probability that some coin will get 10 heads?

**Answer:** $\approx 63\%$

© ⓜ Creator: Yaser Abu-Mostafa - LFD Lecture 2

14/17

# Coin analogy

**1** $\frac{1}{2} \times \frac{1}{2} \times \cdots \times \frac{1}{2}$ (10 times) is $\approx \frac{1}{1000} = 0.1\%$.

**2** The probability of not getting 10 heads for one coin is $\approx (1 - \frac{1}{1000})$.

**3** The probability of not getting 10 heads for any of 1000 coins is $\approx (1 - \frac{1}{1000})^{1000}$.

**4** $\lim_{n \to \infty} (1 - \frac{1}{n})^n = \frac{1}{e}$.

**5** $\frac{1}{e} \approx \frac{1}{2.718} \approx 0.37$

**6** the probability of this not happening, i.e. at least one coin of the 1000 coins will give 10 heads, is 1 minus $0.37 = 0.63$

# From coins to learning



From coins to learning

BINGO?

© Creator: Yaser Abu-Mostafa – LFD Lecture 2

15/17

# From coins to learning

The Hoeffding inequality applies to each bin individually. The inequality states that

$$\mathbb{P}[|E_{\text{in}}(h) - E_{\text{out}}(h)| > \epsilon] \leq 2e^{-2\epsilon^2 N} \quad \text{for any} \ \ \epsilon > 0$$

where

1. the hypothesis $h$ is fixed before the data is generated,

2. the probability is with respect to random data sets $\mathcal{D}$.

The assumption "$h$ is fixed before the data set is generated" is critical to the validity of the bound.

# From coins to learning

In learning, we consider an entire hypothesis set, say $\mathcal{H} = \{h_1, h_2, \ldots, h_M\}$ (with a finite number of hypotheses), instead of just one hypothesis $h$. Then, the learning algorithm picks the final hypothesis $g$ based on $\mathcal{D}$.

The statement we would like to make is **not**

$$\mathbb{P}[|E_{\text{in}}(h_m) - E_{\text{out}}(h_m)| > \epsilon] \text{ is small for any fixed } h_m \in \mathcal{H},$$

where $m = 1, 2, \ldots, M$, but **rather**

$$\mathbb{P}[|E_{\text{in}}(g) - E_{\text{out}}(g)| > \epsilon] \text{ is small for the final hypothesis } g.$$

# A simple solution

$|E_{\text{in}}(g) - E_{\text{out}}(g)| > \epsilon$

$$\implies$$

$$|E_{\text{in}}(h_1) - E_{\text{out}}(h_1)| > \epsilon$$
$$\text{or } |E_{\text{in}}(h_2) - E_{\text{out}}(h_2)| > \epsilon$$

$\cdots$

$$\text{or } |E_{\text{in}}(h_M) - E_{\text{out}}(h_M)| > \epsilon$$

# A simple solution

A simple solution

$$\mathbb{P}[\,|E_{\text{in}}(g) - E_{\text{out}}(g)| > \epsilon\,] \;\leq\; \mathbb{P}[\quad |E_{\text{in}}(h_1) - E_{\text{out}}(h_1)| > \epsilon$$

$$\textbf{or } |E_{\text{in}}(h_2) - E_{\text{out}}(h_2)| > \epsilon$$

$$\cdots$$

$$\textbf{or } |E_{\text{in}}(h_M) - E_{\text{out}}(h_M)| > \epsilon\,]$$

$$\leq\; \sum_{m=1}^{M} \mathbb{P}\left[|E_{\text{in}}(h_m) - E_{\text{out}}(h_m)| > \epsilon\right]$$

# The final verdict

## The final verdict

$$\mathbb{P}[\ |E_{\text{in}}(g) - E_{\text{out}}(g)| > \epsilon\ ] \leq \sum_{m=1}^{M} \mathbb{P}\left[|E_{\text{in}}(h_m) - E_{\text{out}}(h_m)| > \epsilon\right]$$

$$\leq \sum_{m=1}^{M} 2e^{-2\epsilon^2 N}$$

$$\mathbb{P}[|E_{\text{in}}(g) - E_{\text{out}}(g)| > \epsilon] \leq 2Me^{-2\epsilon^2 N}$$

# Generalization error

- The out-of-sample error $E_{out}$ measures how well our training on $\mathcal{D}$ has generalized to unseen data points. $E_{out}$ is based on the performance over the entire input space $\mathcal{X}$.

- The in-sample error $E_{in}$ is based on the training data points.

- The *generalization error* is the discrepancy between $E_{in}$ and $E_{out}$. Generalization error is also used as another name for $E_{out}$ (but not in this unit).

- The Hoeffding inequality provides a way to charaterize the generalization error with a probabilistic bound.

# Generalization bound

The Hoeffding inequality states that

$$\mathbb{P}[|E_{\text{in}}(g) - E_{\text{out}}(g)| > \epsilon] \leq 2Me^{-2\epsilon^2 N} \quad \text{for any } \epsilon > 0$$

This can be rephrased as follows. Pick a tolerance level $\delta$, for example $\delta = 0.01$, and assert with probability at least $1 - \delta$ that

$$E_{\text{out}}(g) \leq E_{\text{in}}(g) + \sqrt{\frac{1}{2N} \ln\left(\frac{2M}{\delta}\right)}.$$

This is called a *generalization bound* since it bounds $E_{\text{out}}$ in terms of $E_{\text{in}}$

# Generalization bound

$$\mathbb{P}[|E_{\text{in}}(g) - E_{\text{out}}(g)| > \epsilon] \leq 2Me^{-2\epsilon^2 N}$$

$$\implies 1 - \mathbb{P}[|E_{\text{in}}(g) - E_{\text{out}}(g)| \leq \epsilon] \leq 2Me^{-2\epsilon^2 N}$$

$$\implies \mathbb{P}[|E_{\text{in}}(g) - E_{\text{out}}(g)| \leq \epsilon] \geq 1 - 2Me^{-2\epsilon^2 N}$$

- With probability at least $1 - 2Me^{-2\epsilon^2 N}$,
  $|E_{\text{in}}(g) - E_{\text{out}}(g)| \leq \epsilon$, which implies $E_{\text{out}}(g) \leq E_{\text{in}}(g) + \epsilon$

- If $\delta = 2Me^{-2\epsilon^2 N}$, then $\epsilon = \sqrt{\frac{1}{2N} \ln\left(\frac{2M}{\delta}\right)}$

# Note

- The error bound depends on $M$, the size of the hypothesis set $\mathcal{H}$

- If $\mathcal{H}$ is an infinite set, the bound goes to infinity and becomes useless

- $M$ can be replaced with something finite (the <u>effective</u> number of hypotheses), so that the bound is meaningful.

$$\mathbb{P}\left[\,|E_{\text{in}}(g) - E_{\text{out}}(g)| > \epsilon\,\right] \leq 4\ m_{\mathcal{H}}(2N)\ e^{-\frac{1}{8}\epsilon^2 N}$$

The Vapnik-Chervonenkis Inequality

**We will not cover the VC Inequality.**

# Outline

# Feasibility of learning

- If we insist on a deterministic answer, i.e. $\mathcal{D}$ tells us something <u>certain</u> about $f$ outside of $\mathcal{D}$, then the answer is no.

- If we accept a probabilistic answer, i.e. $\mathcal{D}$ tells us something <u>likely</u> about $f$ outside of $\mathcal{D}$, then the answer is yes.

# Feasibility of learning

What we know so far

Learning is feasible. It is likely that

$$E_{\text{out}}(g) \approx E_{\text{in}}(g)$$

Is this learning?

We need $g \approx f$, which means

$$E_{\text{out}}(g) \approx 0$$

# Feasibility of learning

The 2 questions of learning

$E_{out}(g) \approx 0$ is achieved through:

$$\underbrace{E_{out}(g) \approx E_{in}(g)}_{\text{Lecture 2}} \qquad \text{and} \qquad \underbrace{E_{in}(g) \approx 0}_{\text{Lecture 3}}$$

Learning is thus split into 2 questions:

1. Can we make sure that $E_{out}(g)$ is close enough to $E_{in}(g)$?

2. Can we make $E_{in}(g)$ small enough?

21/22

- The Hoeffding Inequality addresses question 1 only.

- We answer the second question after running the learning algorithm on the the training data.

# Feasibility of learning

**The complexity of** $\mathcal{H}$.

Question 1: According to the Hoeffding Inequality, a larger $M$ increases the risk that $E_{\text{in}}(g)$ will be a poor estimate of $E_{\text{out}}(g)$ $\implies$ we need to control $M$ (a measure of the complexity of $\mathcal{H}$).

Question 2: We stand a better chance if $\mathcal{H}$ is more complex $\implies$ a more complex $\mathcal{H}$ gives us more flexibility in finding some $g$ that fits the data well.

**The complexity of** $f$.

Question 1: If we fix the hypothesis set and the number of training examples, the inequality provides the same bound $\implies$ The complexity of $f$ does not affect how well $E_{\text{in}}(g)$ approximates $E_{\text{out}}(g)$.

Question 2: The data from a complex $f$ are harder to fit than the data from a simple $f$ (large $E_{\text{in}}(g)$ ). We can increase the complexity of $\mathcal{H}$, but then $E_{\text{out}}(g)$ will not be as close to $E_{\text{in}}(g)$.

# Approximation-generalization tradeoff

## Approximation-generalization tradeoff

Small $E_{\text{out}}$: good approximation of $f$ out of sample.

More complex $\mathcal{H} \implies$ better chance of **approximating** $f$

Less complex $\mathcal{H} \implies$ better chance of **generalizing** out of sample

Ideal $\mathcal{H} = \{f\}$      winning lottery ticket ☺

# Quantifying the tradeoff

VC analysis was one approach: $E_{\text{out}} \leq E_{\text{in}} + \Omega$

Bias-variance analysis is another: decomposing $E_{\text{out}}$ into

    **1.** How well $\mathcal{H}$ can approximate $f$

    **2.** How well we can zoom in on a good $h \in \mathcal{H}$

Applies to **real-valued targets** and uses **squared error**

# Bias and variance decomposition

Start with $E_{\text{out}}$

$$E_{\text{out}}(g^{(\mathcal{D})}) = \mathbb{E}_{\mathbf{x}}\left[\left(g^{(\mathcal{D})}(\mathbf{x}) - f(\mathbf{x})\right)^2\right]$$

$$\mathbb{E}_{\mathcal{D}}\left[E_{\text{out}}(g^{(\mathcal{D})})\right] = \mathbb{E}_{\mathcal{D}}\left[\mathbb{E}_{\mathbf{x}}\left[\left(g^{(\mathcal{D})}(\mathbf{x}) - f(\mathbf{x})\right)^2\right]\right]$$

$$= \mathbb{E}_{\mathbf{x}}\left[\mathbb{E}_{\mathcal{D}}\left[\left(g^{(\mathcal{D})}(\mathbf{x}) - f(\mathbf{x})\right)^2\right]\right]$$

Now, let us focus on:

$$\mathbb{E}_{\mathcal{D}}\left[\left(g^{(\mathcal{D})}(\mathbf{x}) - f(\mathbf{x})\right)^2\right]$$

# Bias and variance decomposition

Bias and variance

$$\mathbb{E}_{\mathcal{D}}\left[\left(g^{(\mathcal{D})}(\mathbf{x}) - f(\mathbf{x})\right)^2\right] = \underbrace{\mathbb{E}_{\mathcal{D}}\left[\left(g^{(\mathcal{D})}(\mathbf{x}) - \bar{g}(\mathbf{x})\right)^2\right]}_{\text{var}(\mathbf{x})} + \underbrace{\left(\bar{g}(\mathbf{x}) - f(\mathbf{x})\right)^2}_{\text{bias}(\mathbf{x})}$$

Therefore, $\mathbb{E}_{\mathcal{D}}\left[E_{\text{out}}(g^{(\mathcal{D})})\right] = \mathbb{E}_{\mathbf{x}}\left[\mathbb{E}_{\mathcal{D}}\left[\left(g^{(\mathcal{D})}(\mathbf{x}) - f(\mathbf{x})\right)^2\right]\right]$

$$= \mathbb{E}_{\mathbf{x}}[\text{bias}(\mathbf{x}) + \text{var}(\mathbf{x})]$$

$$= \quad \text{bias} \quad + \quad \text{var}$$

8/22

# Outline

**1** Verification vs learning

**2** Feasibility of learning

**3** Error measures

**4** Noisy targets

# Learning diagram



The learning diagram – where we left it

**UNKNOWN TARGET FUNCTION**
$f: X \to Y$

**PROBABILITY DISTRIBUTION**
$P$ on $X$

$x_1, \dots, x_N$

**TRAINING EXAMPLES**
$(x_1, y_1), \dots, (x_N, y_N)$

**LEARNING ALGORITHM**
$A$

**FINAL HYPOTHESIS**
$g: X \to Y$

**HYPOTHESIS SET**
$H$

# Error measures

What does "$h \approx f$" mean?

Error measure: $E(h, f)$

Almost always *pointwise definition*: $\mathrm{e}\big(h(\mathbf{x}), f(\mathbf{x})\big)$

Examples:

Squared error: $\quad \mathrm{e}\big(h(\mathbf{x}), f(\mathbf{x})\big) = \big(h(\mathbf{x}) - f(\mathbf{x})\big)^2$

Binary error: $\quad \mathrm{e}\big(h(\mathbf{x}), f(\mathbf{x})\big) = [\![h(\mathbf{x}) \neq f(\mathbf{x})]\!]$

# From pointwise to overall

## From pointwise to overall

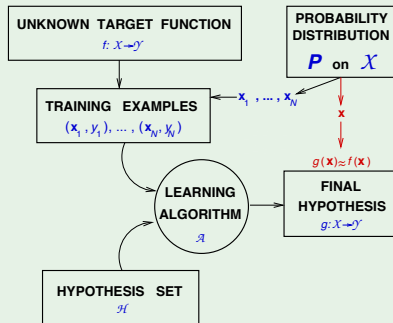Overall error $E(h, f)$ = average of pointwise errors $\mathrm{e}\left(h(\mathbf{x}), f(\mathbf{x})\right)$.

In-sample error:

$$E_{\mathrm{in}}(h) = \frac{1}{N} \sum_{n=1}^{N} \mathrm{e}\left(h(\mathbf{x}_n), f(\mathbf{x}_n)\right)$$

Out-of-sample error:

$$E_{\mathrm{out}}(h) = \mathbb{E}_{\mathbf{x}}\left[\mathrm{e}\left(h(\mathbf{x}), f(\mathbf{x})\right)\right]$$

# Learning diagram updated



The learning diagram – with pointwise error

UNKNOWN TARGET FUNCTION
$f: \mathcal{X} \to \mathcal{Y}$

PROBABILITY DISTRIBUTION
$P$ on $\mathcal{X}$

$x_1, \ldots, x_N$

TRAINING EXAMPLES
$(x_1, y_1), \ldots, (x_N, y_N)$

$x$

$g(x) \approx f(x)$

LEARNING ALGORITHM
$\mathcal{A}$

FINAL HYPOTHESIS
$g: \mathcal{X} \to \mathcal{Y}$

HYPOTHESIS SET
$\mathcal{H}$

9/22

# How to choose the error measure?

How to choose the error measure

Fingerprint verification:

Two types of error:

*false accept*  and  *false reject*

How do we penalize each type?



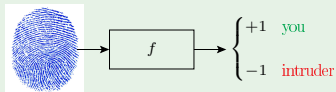|   |   | $f$ | |
|---|---|---|---|
|   |   | $+1$ | $-1$ |
| $h$ | $+1$ | no error | *false accept* |
|   | $-1$ | *false reject* | no error |

# The supermarket example

## The error measure - for supermarkets

Supermarket verifies fingerprint for discounts

False reject is costly; customer gets annoyed!

False accept is minor; gave away a discount
and intruder left their fingerprint ☺



$$\begin{array}{cc|cc} & & \multicolumn{2}{c}{f} \\ & & +1 & -1 \\ \hline h & +1 & 0 & 1 \\ & -1 & 10 & 0 \end{array}$$
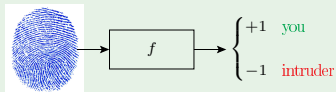
# The CIA example

The error measure -    for the CIA

CIA verifies fingerprint for security
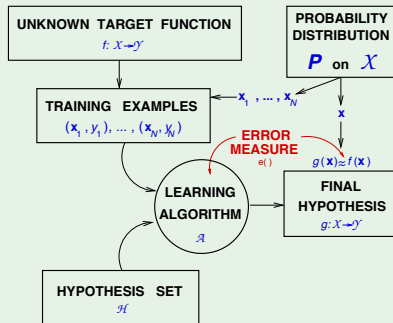
False accept is a disaster!

False reject can be tolerated
Try again; you are an employee ☺



$$\begin{array}{c|cc} & \multicolumn{2}{c}{f} \\ & +1 & -1 \\ \hline h \quad +1 & 0 & 1000 \\ -1 & 1 & 0 \end{array}$$

# Learning diagram updated

The learning diagram - with error measure

# Outline

# Noisy targets

The 'target function' is not always a *function*

Consider the credit-card approval:

| age | 23 years |
|---|---|
| annual salary | $30,000 |
| years in residence | 1 year |
| years in job | 1 year |
| current debt | $15,000 |
| . . . | . . . |

two 'identical' customers $\longrightarrow$ two different behaviors

# Target distribution

**Target 'distribution'**

Instead of $y = f(\mathbf{x})$, we use target *distribution*:
$$P(y \mid \mathbf{x})$$

$(\mathbf{x}, y)$ is now generated by the joint distribution:
$$P(\mathbf{x})P(y \mid \mathbf{x})$$

Noisy target = deterministic target $f(\mathbf{x}) = \mathbb{E}(y|\mathbf{x})$ plus noise $y - f(\mathbf{x})$

Deterministic target is a special case of noisy target:
$$P(y \mid \mathbf{x}) \text{ is zero except for } y = f(\mathbf{x})$$

# Final learning diagram



The learning diagram – including noisy target