

CSE4088 Introduction to Machine Learning

Error and Noise

Slides are adopted from lecture notes of Yaser Abu-Mostafa

Review of Last Lecture

- Linear models use the ‘**signal**’:

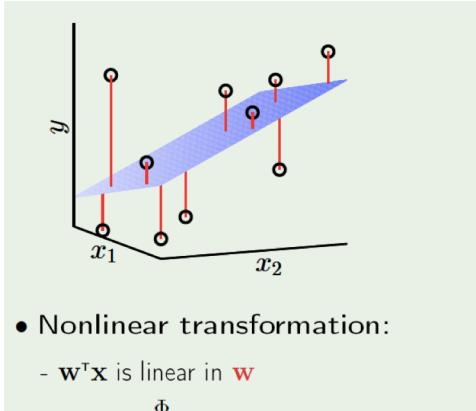
$$\sum_{i=0}^d w_i x_i = \mathbf{w}^\top \mathbf{x}$$

- Classification: $h(\mathbf{x}) = \text{sign}(\mathbf{w}^\top \mathbf{x})$
- Regression: $h(\mathbf{x}) = \mathbf{w}^\top \mathbf{x}$

- Linear regression algorithm:

$$\mathbf{w} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

“one-step learning”

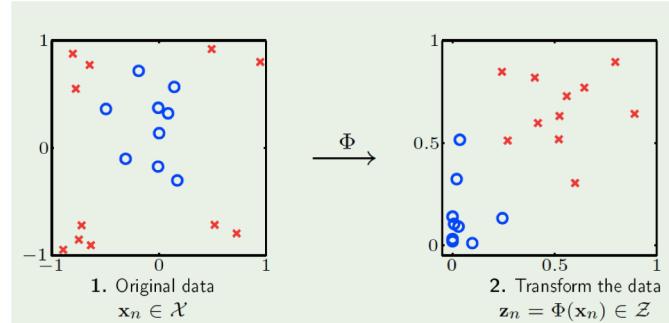


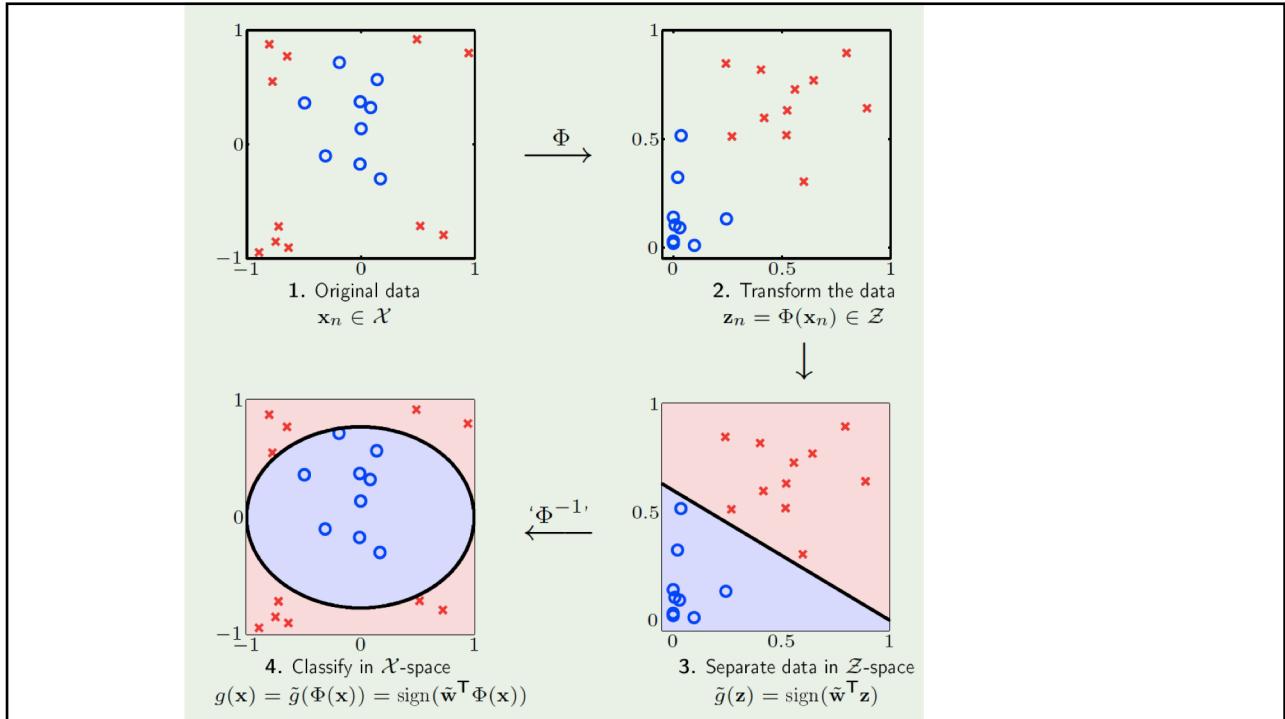
- Nonlinear transformation:

- $\mathbf{w}^\top \mathbf{x}$ is linear in \mathbf{w}
- Any $\mathbf{x} \xrightarrow{\Phi} \mathbf{z}$ preserves this linearity.
- Example: $(x_1, x_2) \xrightarrow{\Phi} (x_1^2, x_2^2)$

Outline

- Nonlinear transformation (continued)
- Error measures
- Noisy targets
- Preamble to the theory





What transforms to what

$$\mathbf{x} = (x_0, x_1, \dots, x_d) \xrightarrow{\Phi} \mathbf{z} = (z_0, z_1, \dots, z_{\tilde{d}})$$

$$\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N \xrightarrow{\Phi} \mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_N$$

$$y_1, y_2, \dots, y_N \xrightarrow{\Phi} y_1, y_2, \dots, y_N$$

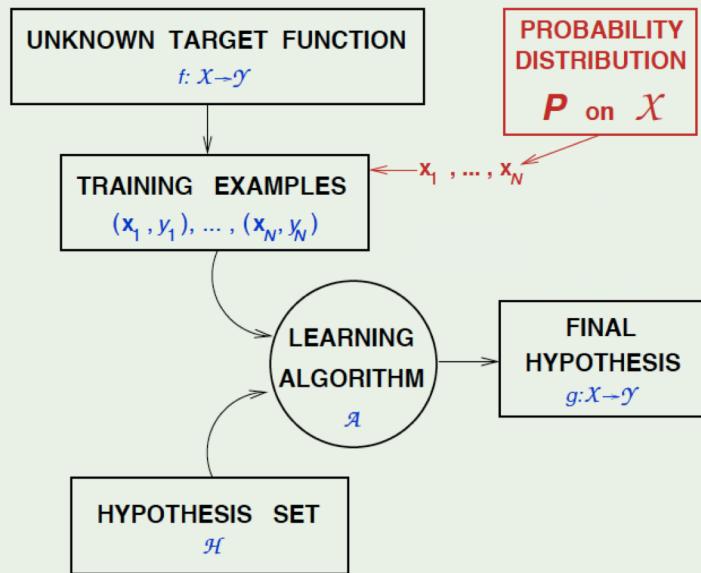
No weights in \mathcal{X} $\tilde{\mathbf{w}} = (w_0, w_1, \dots, w_{\tilde{d}})$

$$g(\mathbf{x}) = \text{sign}(\tilde{\mathbf{w}}^\top \Phi(\mathbf{x}))$$

Outline

- Nonlinear transformation (continued)
- Error measures
- Noisy targets
- Preamble to the theory

The learning diagram - where we left it



Error Measures

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

Error measure: $E(h, f)$

Error Measures

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

Error measure: $E(h, f)$

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

Error Measures

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

Error measure: $E(h, f)$

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

Examples:

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

$$\text{Binary error: } e(h(\mathbf{x}), f(\mathbf{x})) = [h(\mathbf{x}) \neq f(\mathbf{x})]$$

From pointwise to overall

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

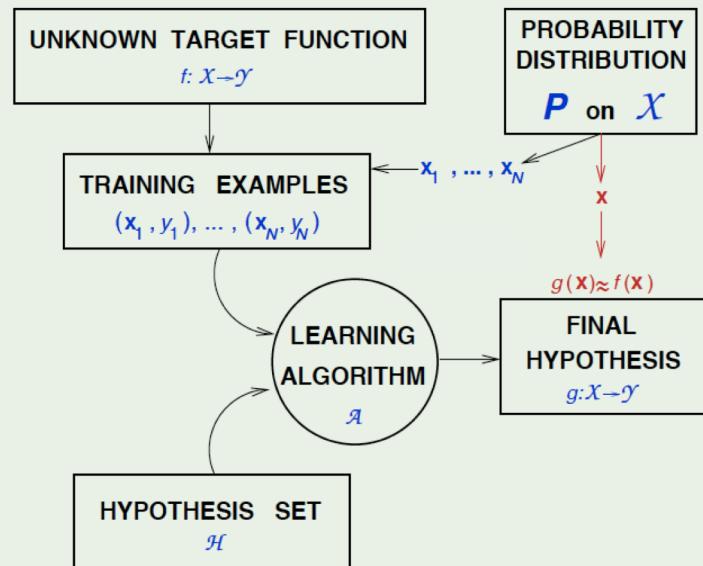
In-sample error:

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

Out-of-sample error:

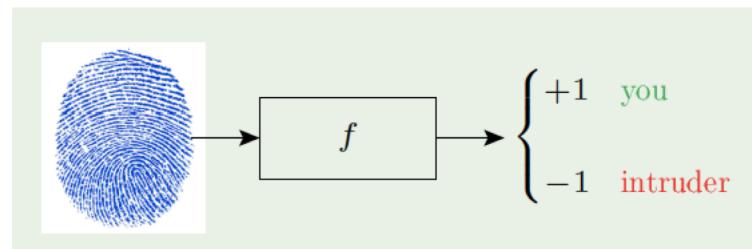
$$E_{\text{out}}(h) = \mathbb{E}_{\mathbf{x}}[e(h(\mathbf{x}), f(\mathbf{x}))]$$

The learning diagram - with pointwise error



How to choose the error measure

Fingerprint verification:



Two types of error:

false accept and *false reject*

How to choose the error measure

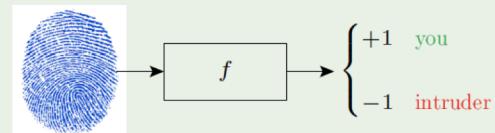
Fingerprint verification:

Two types of error:

false accept and *false reject*

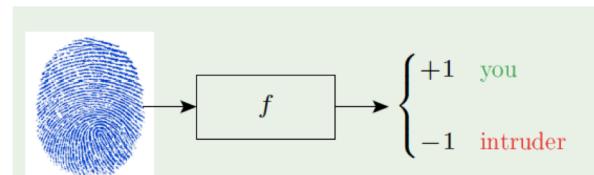
How do we penalize each type?

		<i>f</i>	
		+1	-1
<i>h</i>	+1	no error	<i>false accept</i>
	-1	<i>false reject</i>	no error



The error measure – for supermarkets

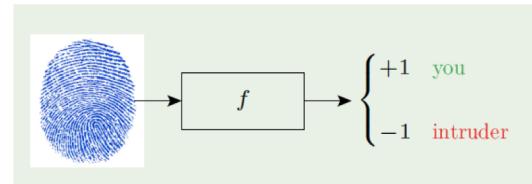
- Supermarket verifies fingerprint for discounts.
- False reject is costly; customer gets annoyed!
- False accept has minor cost; gave away a discount and intruder left their fingerprint ☺



		<i>f</i>	
		+1	-1
<i>h</i>	+1	0	1
	-1	10	0

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 ☺

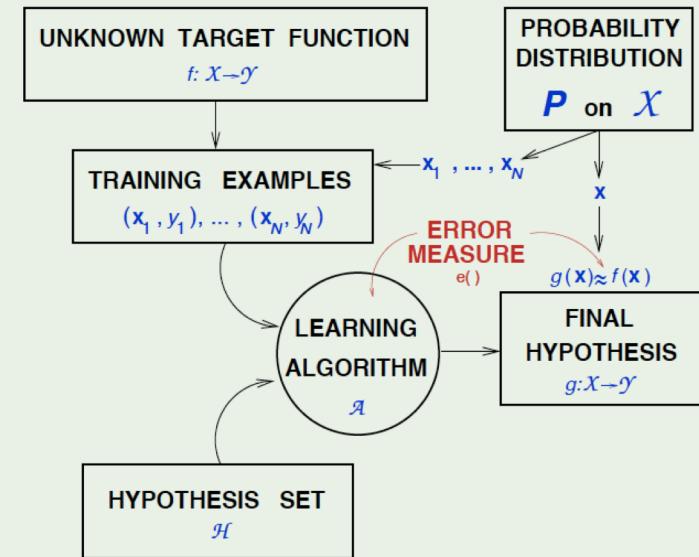


		f	
		+1	-1
h	+1	0	1000
	-1	1	0

Take-home lesson

- The error measure should be specified by the user.
- Not always possible. Alternatives:
 - **Plausible** measures: squared error – Gaussian noise
 - **Friendly** measures: closed-form solution (as in regression), convex optimization

The learning diagram - with error measure



Outline

- Nonlinear transformation (continued)
- Error measures
- Noisy targets
- Preamble to the theory

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’ customer → two different behaviors

Target ‘distribution’

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

$$P(y | \mathbf{x})$$

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

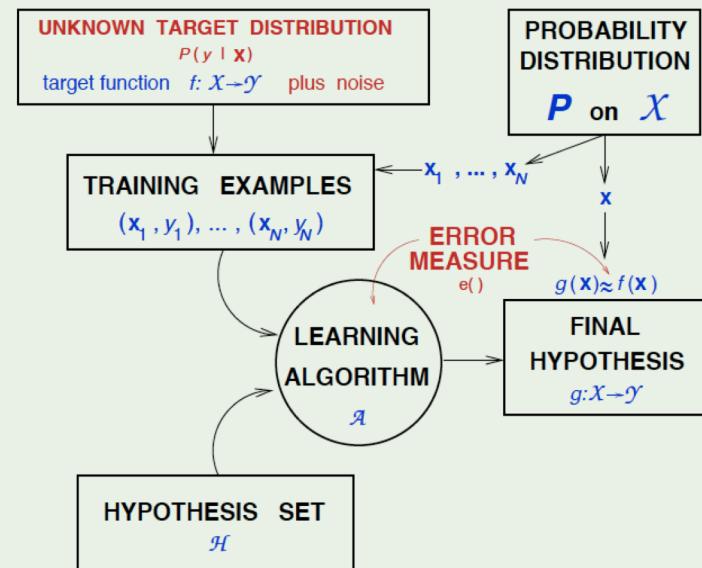
$$P(\mathbf{x})P(y | \mathbf{x})$$

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:

$$P(y | \mathbf{x}) \text{ is zero except for } y = f(\mathbf{x})$$

The learning diagram - including noisy target



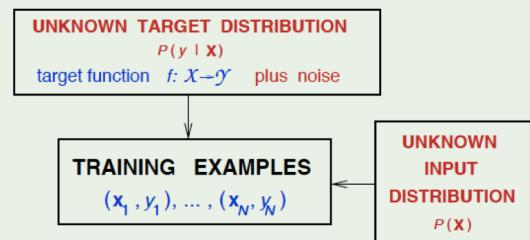
Distinction between $P(y|x)$ and $P(x)$

Both convey probabilistic aspects of x and y

The target distribution $P(y | x)$
is what we are trying to learn

The input distribution $P(x)$
quantifies relative importance of x

Merging $P(x)P(y|x)$ as $P(x, y)$
mixes the two concepts



Outline

- Nonlinear transformation (continued)
- Error measures
- Noisy targets
- Preamble to the theory

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

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

The 2 questions of learning

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

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

Learning is thus split into 2 questions:

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

What the theory will achieve

Characterizing the feasibility of learning for infinite M

Characterizing the tradeoff:

Model complexity	\uparrow	E_{in}	\downarrow
Model complexity	\uparrow	$E_{\text{out}} - E_{\text{in}}$	\uparrow

