

## Review of Lecture 3

- Linear models use the ‘**signal**’:

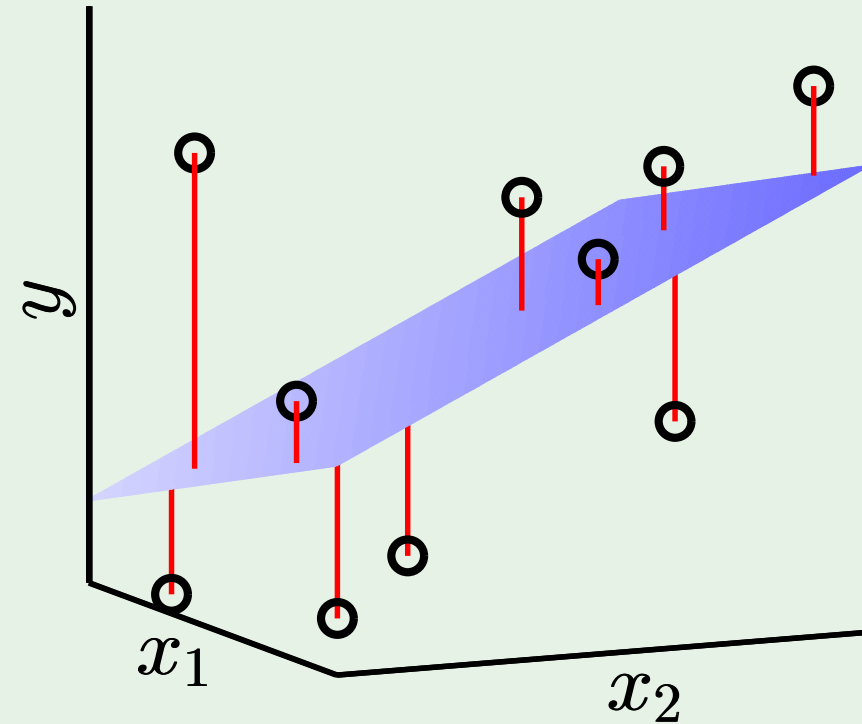
$$\sum_{i=0}^d w_i x_i = \mathbf{w}^\top \mathbf{x}$$

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

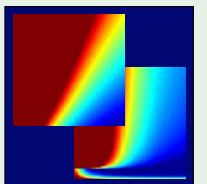
# Learning From Data

Yaser S. Abu-Mostafa  
*California Institute of Technology*

## Lecture 4: **Error and Noise**

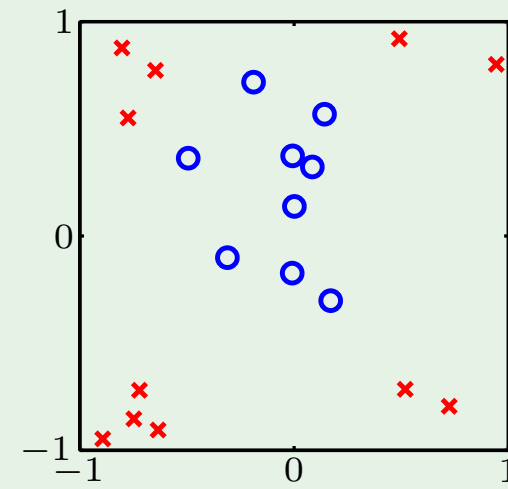


Sponsored by Caltech's Provost Office, E&AS Division, and IST • Thursday, April 12, 2012



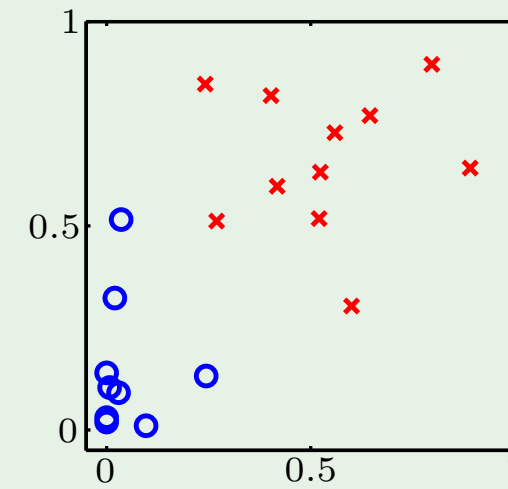
# Outline

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



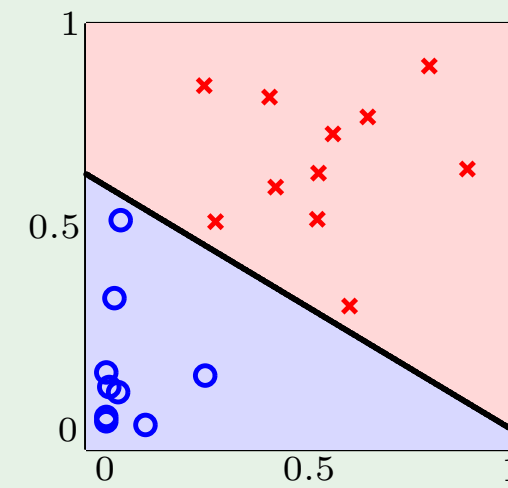
1. Original data  
 $\mathbf{x}_n \in \mathcal{X}$

$\xrightarrow{\Phi}$



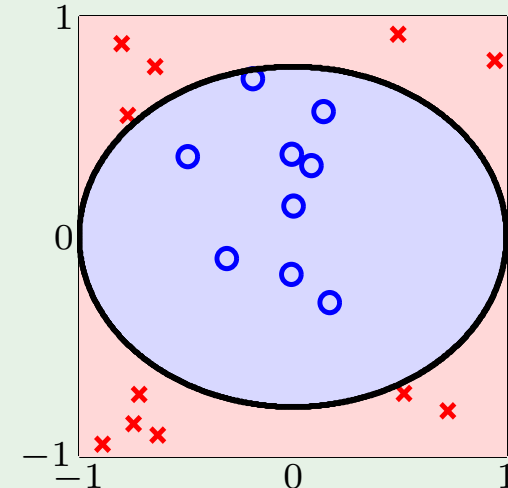
2. Transform the data  
 $\mathbf{z}_n = \Phi(\mathbf{x}_n) \in \mathcal{Z}$

$\downarrow$



3. Separate data in  $\mathcal{Z}$ -space  
 $\tilde{g}(\mathbf{z}) = \text{sign}(\tilde{\mathbf{w}}^T \mathbf{z})$

$\xleftarrow{\Phi^{-1}}$



4. Classify in  $\mathcal{X}$ -space  
 $g(\mathbf{x}) = \tilde{g}(\Phi(\mathbf{x})) = \text{sign}(\tilde{\mathbf{w}}^T \Phi(\mathbf{x}))$

Can transform from 2d (or however many you data has) to however many you require to separate - note however that the wrong transformation can lead to poor generalization

# What transforms to what

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

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

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

e.g. in classification, each training example has the same class in Z space

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

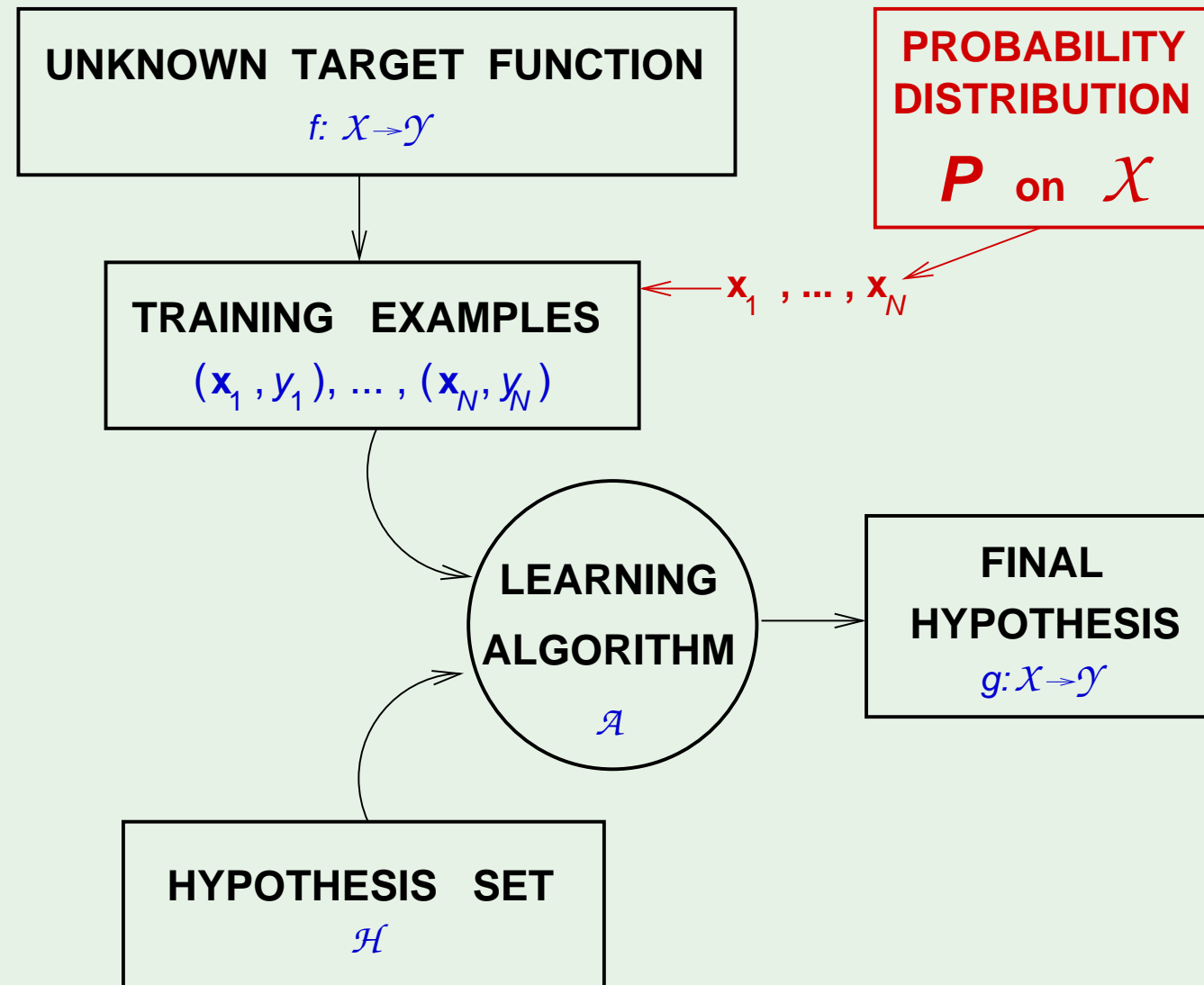
$$g(\mathbf{x}) = \text{sign}(\tilde{\mathbf{w}}^\top \Phi(\mathbf{x}))$$

As above, since the y in X and Z are the same (no transformation), points are classified the same in X and Z space

# 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:  $E(h, f)$  technically a functional, rather than a function, since it returns a number based on two functions, rather than variables - it quantifies how far our hypothesis is from the target.

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

Examples:

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

Binary error: 
$$e(h(\mathbf{x}), f(\mathbf{x})) = \llbracket h(\mathbf{x}) \neq f(\mathbf{x}) \rrbracket$$
 Again,  $\llbracket \dots \rrbracket$  returns 1 if true, 0 if false



## From pointwise to overall

Overall error  $E(h, f)$  = 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:

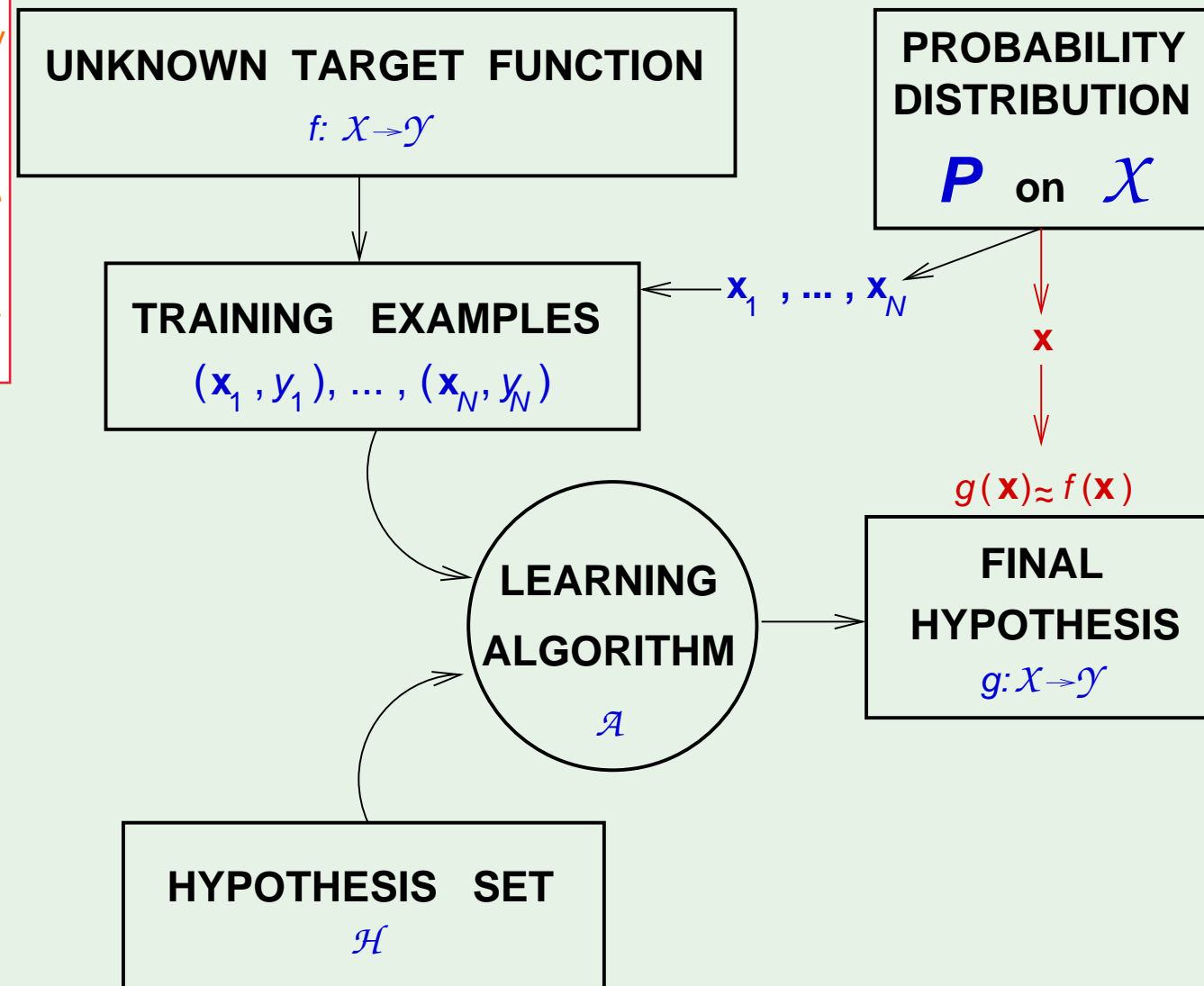
represents error over the whole space (hence expectation value)

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

expectation value with respect to  $\mathbf{x}$ , where  $\mathbf{x}$  is a general point in space - e.g. binary error in this formula gives the probability of error overall

# The learning diagram - with pointwise error

The choice of an error measure affects the outcome of the learning process. Different error measures may lead to different choices of the final hypothesis, even if the target and the data are the same, since the value of a particular error measure may be small while the value of another error measure in the same situation is large. One may view  $E(h, J)$  as the 'cost' of using  $h$  when you should use  $f$ . This cost depends on what  $h$  is used for, and cannot be dictated just by our learning techniques.



So when you test the system you trained with a certain probability distribution (for the training data), you must test with points drawn from the same probability distribution (required to invoke Hoeffding (or the counterpart of Hoeffding for more elaborate functions) - so training and test data must be drawn from same  $P$ ).

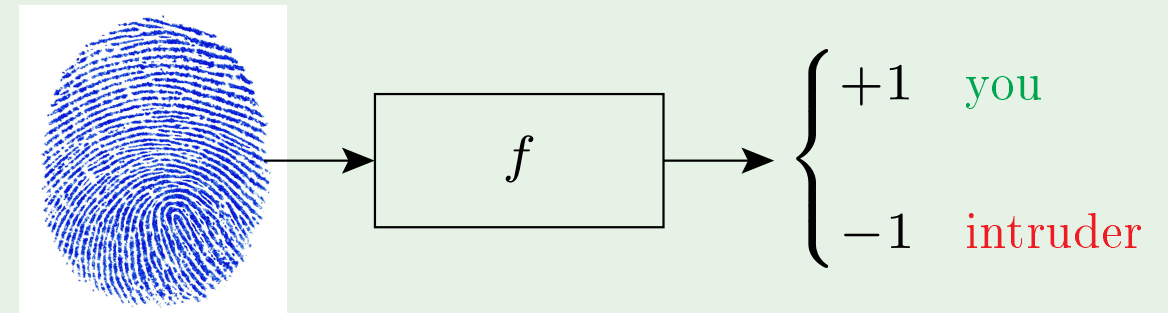
# How to choose the error measure

Fingerprint verification:

Two types of error:

*false accept* and *false reject*

How do we penalize each type?



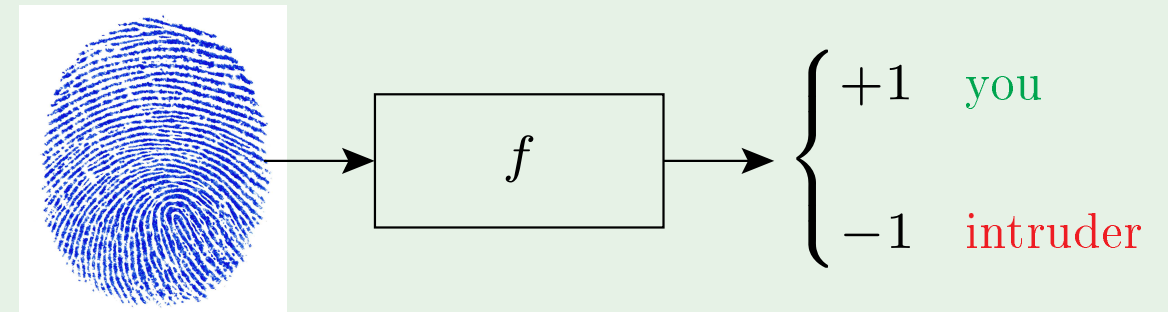
		$f$	
		$+1$	$-1$
$h$	$+1$	no error	<i>false accept</i>
	$-1$	<i>false reject</i>	no error

# The error measure – for supermarkets

Supermarket verifies fingerprint for discounts

False reject is costly; customer gets annoyed!

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



		$f$	
		$+1$	$-1$
$h$	$+1$	0	1
	$-1$	10	0

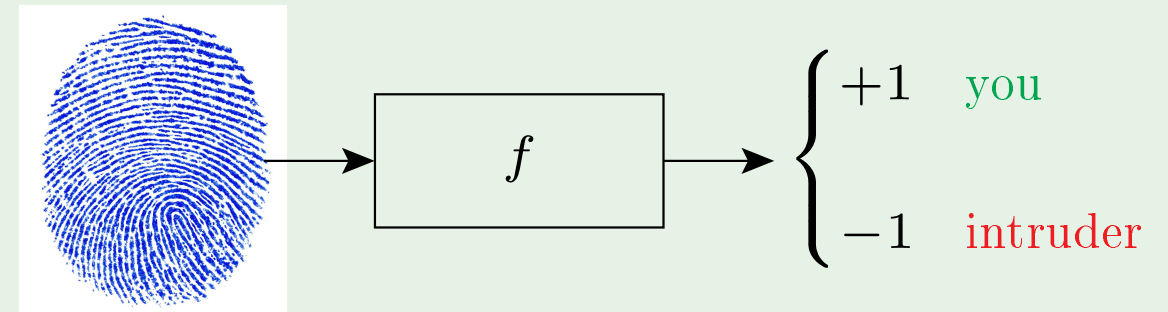
## The error measure – for the CIA

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

From the previous two examples, we see that the error measure is different between two application domains for exactly the same machine learning system (same training data, same target function) - the actual values involved could be determined by assessing the cost of a false accept and false reject to the application domain/user.

# Take-home lesson

The error measure should be specified by the user.

However, this ideal choice may not be possible in practice for two reasons. One is that the user may not provide an error specification, which is not uncommon. The other is that the weighted cost may be a difficult objective function for optimizers to work with.

Not always possible. Alternatives:

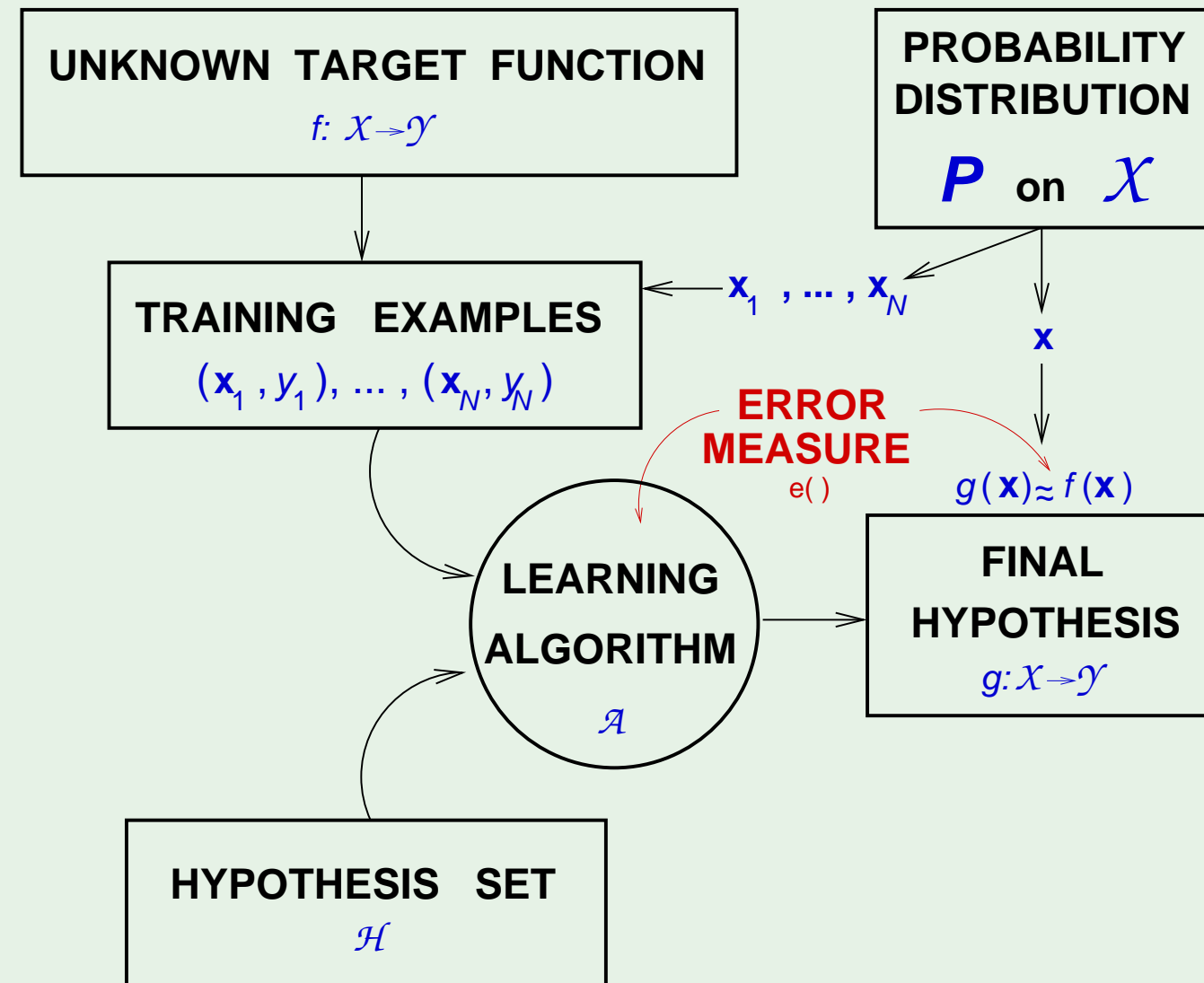
*Plausible* measures: squared error  $\equiv$  Gaussian noise

based on conceptual appeal/merit - e.g. squared error is preferred in general over absolute difference since squared error is optimized over a smooth parabola (has preferable properties) while absolute error has a vertex/discontinuity so becomes a combinatorial optimization instead of a smooth function). Note that if the user specifies that you must minimize the absolute function, it can be worked with. If you are making an analytic choice though, you should pick the friendlier option (in terms of the concept or the optimization).

*Friendly* measures: closed-form solution, convex optimization

practical appeal/easy to use - e.g. we used the least-squared error measure in linear regression, which allowed us to derive a closed-form solution for the calculating the weights in the one-step linear regression learning algorithm, similarly those that allow convex optimization techniques in the determination of weights can be relatively easily solved.

## The learning diagram - with error measure



# 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
...	...

can end up exhibiting

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

- so essentially we have one data point mapping to different outputs, so not a deterministic function (def: returns a unique value for all points in the domain) but a noisy one



# Target 'distribution'

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

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

- some  $y$ 's are more likely for a given  $x$ , instead of there being only one  $y$  for a given  $x$  (in a function) -  $y$  is a random variable that is affected, rather than determined, by the input  $x$ .

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

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

A data point  $(x,y)$  can be treated as a pair generated by the joint distribution,  $P(x)P(y \mid x)$  (assuming independence between the two).

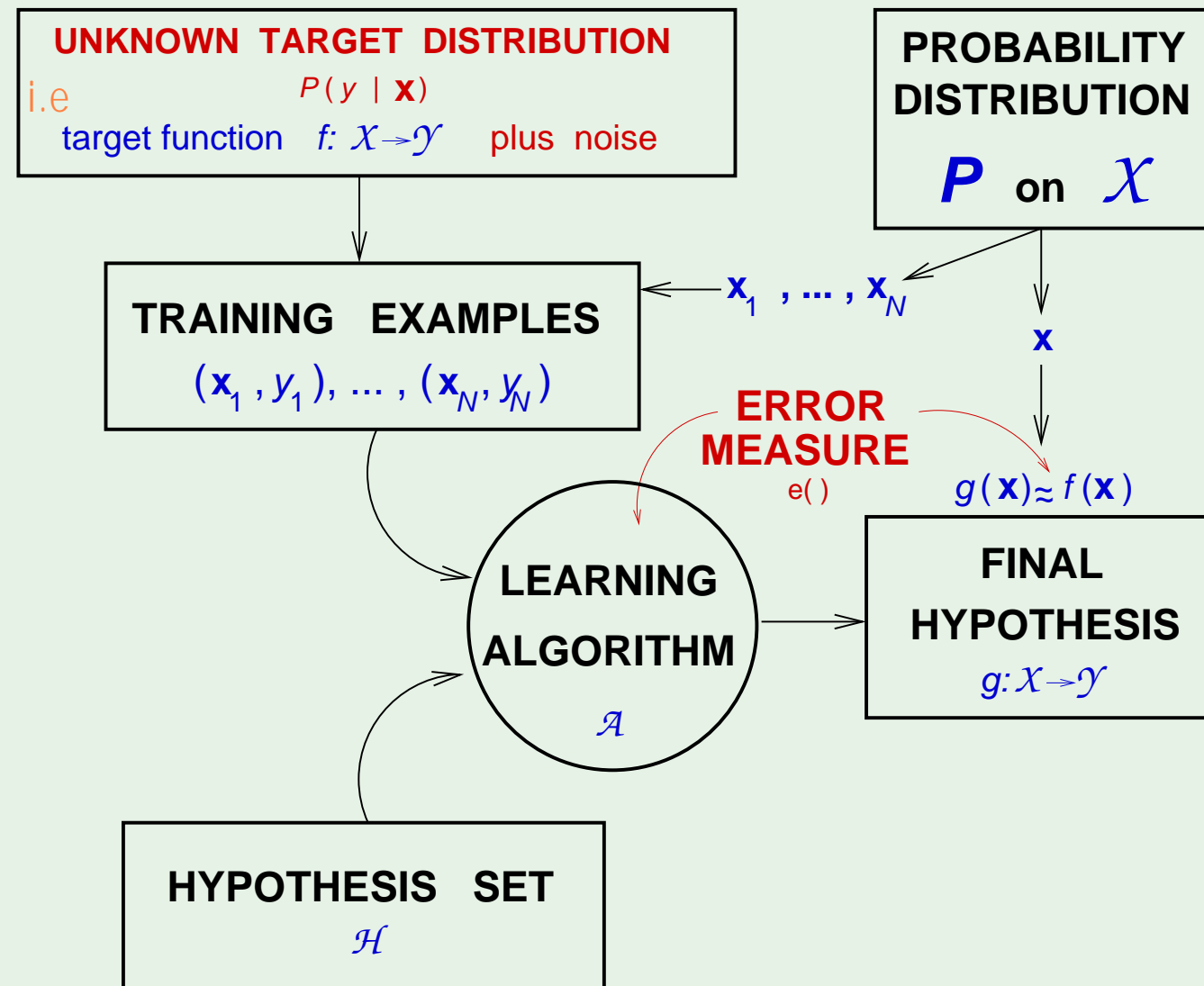
Noisy target = deterministic target  $f(\mathbf{x}) = \mathbb{E}(y|\mathbf{x})$  plus noise  $y - f(\mathbf{x})$

pure noise, so the average should be around zero

Deterministic target is a special case of noisy target:

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

# The learning diagram - including noisy target



# Distinction between $P(y|\mathbf{x})$ and $P(\mathbf{x})$

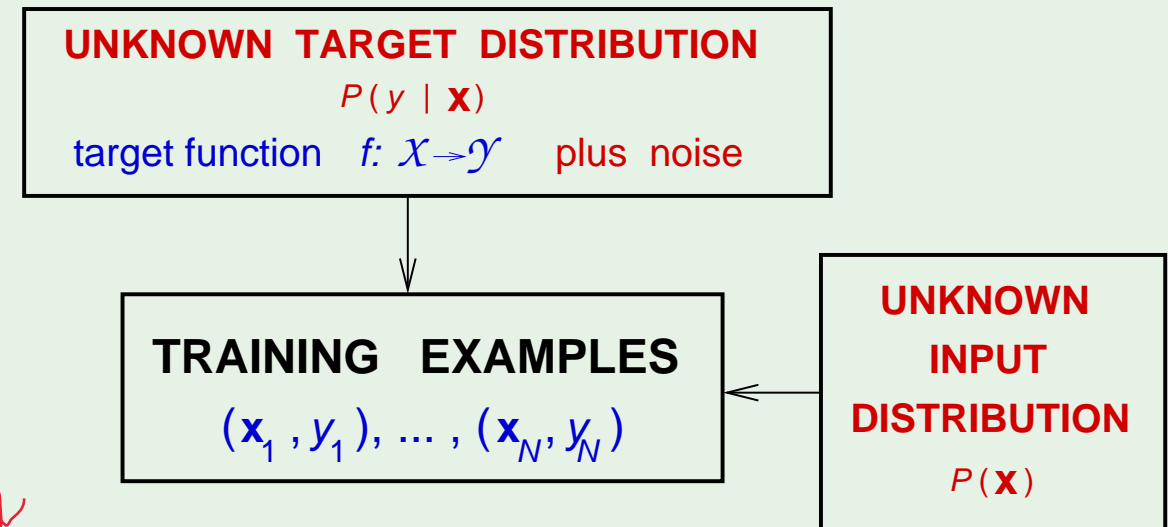
Both convey probabilistic aspects of  $\mathbf{x}$  and  $y$

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

The input distribution  $P(\mathbf{x})$   
quantifies relative importance of  $\mathbf{x}$

Merging  $P(\mathbf{x})P(y|\mathbf{x})$  as  $P(\mathbf{x}, y)$   
mixes the two concepts

target used to generate  
training+test samples



$P(\mathbf{x})$  only quantifies the relative importance of the point  $\mathbf{x}$  in gauging how well we have learned. In the credit approval example: the target distribution is the probability of credit worthiness given the input (e.g. salary) - e.g. decide on the risk of defaulting and assign output +1 (approve credit) with probability 0.9 and reject with 0.1 - this is what we are trying to learn. The input distribution only tells us the distribution of salaries in the general population. So in the instance that  $P(\mathbf{x})$  is skewed to the right (higher salaries) and high salaries lead to good credit worthiness. This model will be tested largely in the comfortable region of high salaries, so a lot of approved credit with little error. However, if the mass of probability is around the borderline cases, the system will perform worse categorizing these borderline cases. So  $P(\mathbf{x})$  does give the weights that will finally grade the hypothesis but we are not trying to learn  $P(\mathbf{x})$

# Outline

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

# What we know so far

by relying on the fact that  
Learning is feasible. It is likely that

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

(which represents generalization -  $E_{\text{in}}$  is a proxy for  $E_{\text{out}}$  which we do not know)

Is this learning?

We need  $g \approx f$ , which means

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

(which means we learned well)

# The 2 questions of learning

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

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

Lecture 2

by Hoeffding

and

$$\underbrace{E_{\text{in}}(g) \approx 0}$$

Lecture 3

We know  $E_{\text{in}}$  and can use a learning algorithm (like linear regression) to make it as small as possible -  $E_{\text{in}} \sim 0$  means we need to learn the training data well (hopefully not overfitting - relevant to  $E_{\text{out}} \sim E_{\text{in}}$ )

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?

Theoretical question - over next 4 lectures

Practical question - over 8 lectures after Q1

# What the theory will achieve

Characterizing the feasibility of learning for  
infinite  $M$

Characterizing the tradeoff:

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

