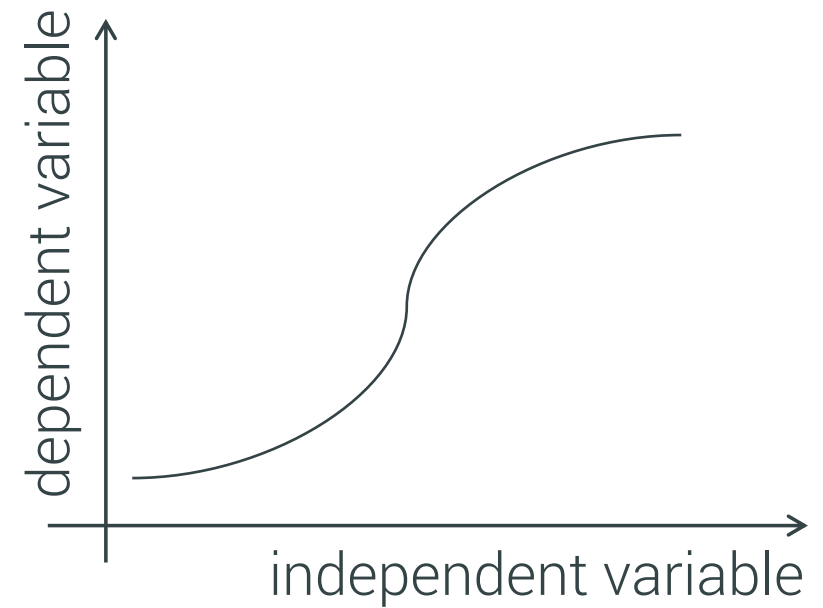


# End-to-end machine learning

## Lecture 02

# Quiz

# Common language



**independent** variable

input

predictor

feature

$x$

**dependent** variable

output

response

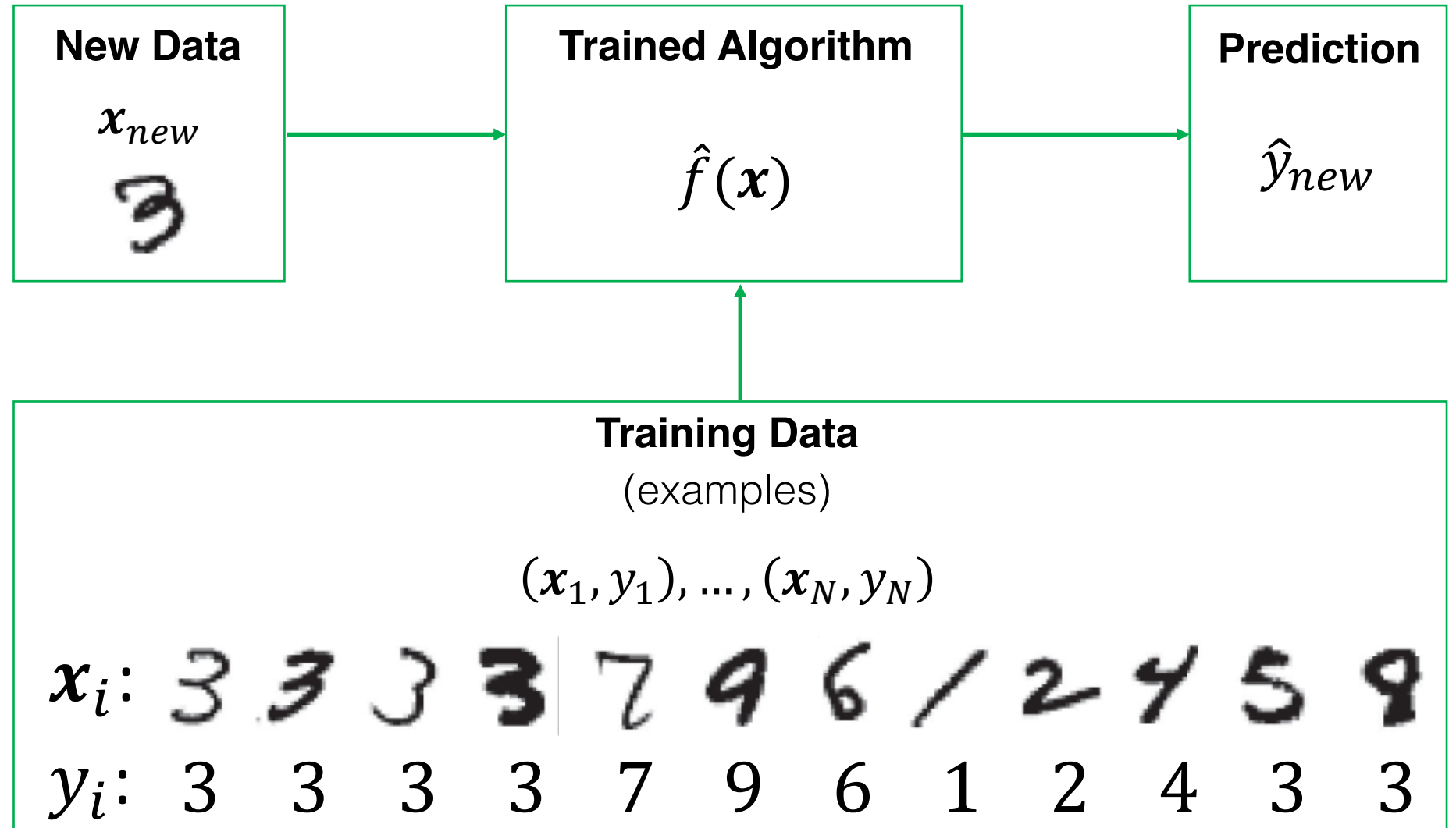
target

$y$

# Supervised learning

Objective: create an algorithm that predicts well

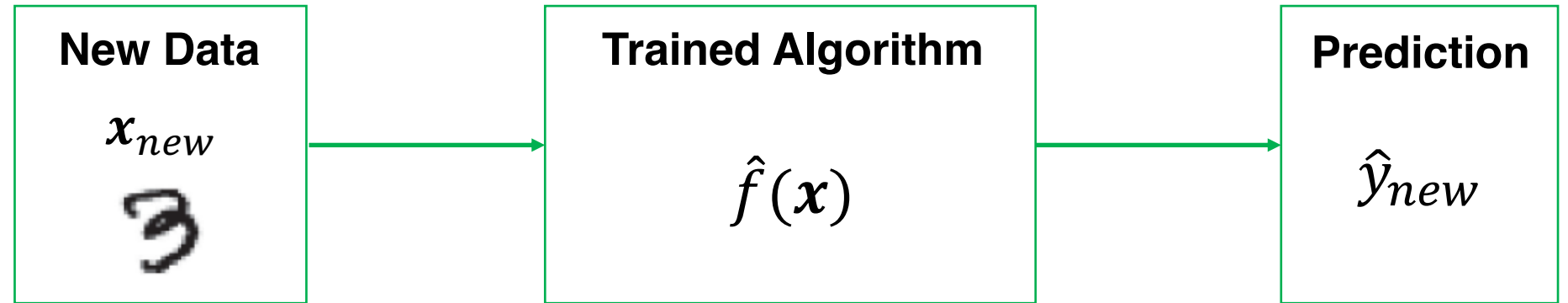
Example:  
**Digits classification**



# Supervised learning

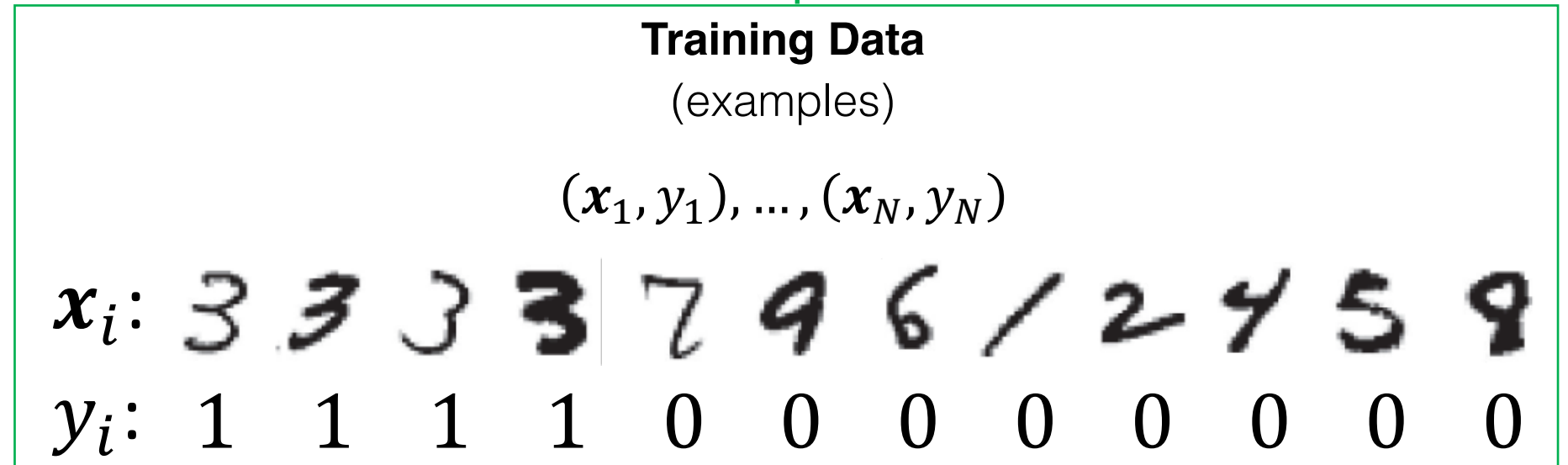
Objective: create an algorithm that predicts well

Example:  
**Digits classification**



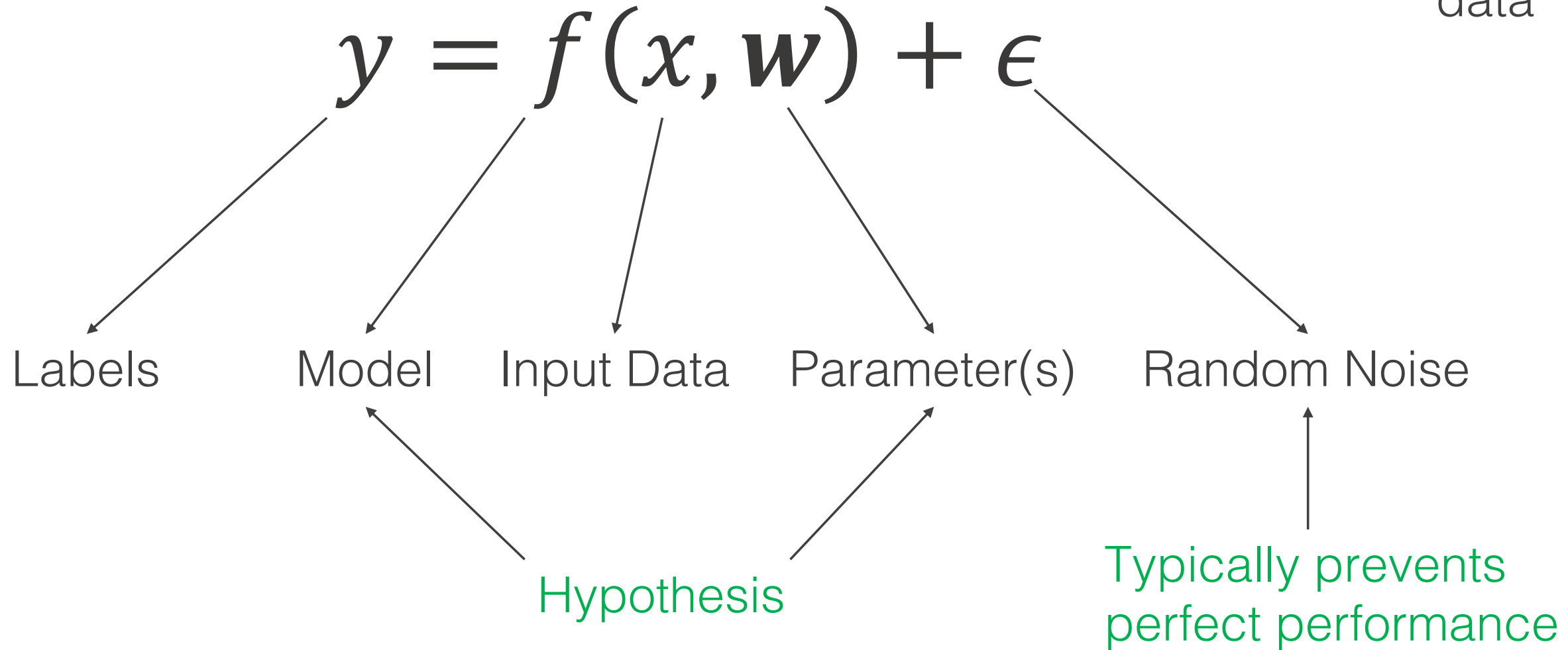
Simplification for example:  
Is the digit a “3”?

This becomes a binary  
classification problem:  
 $y = 1$  implies a “3”  
 $y = 0$  implies another digit



# Supervised machine learning model

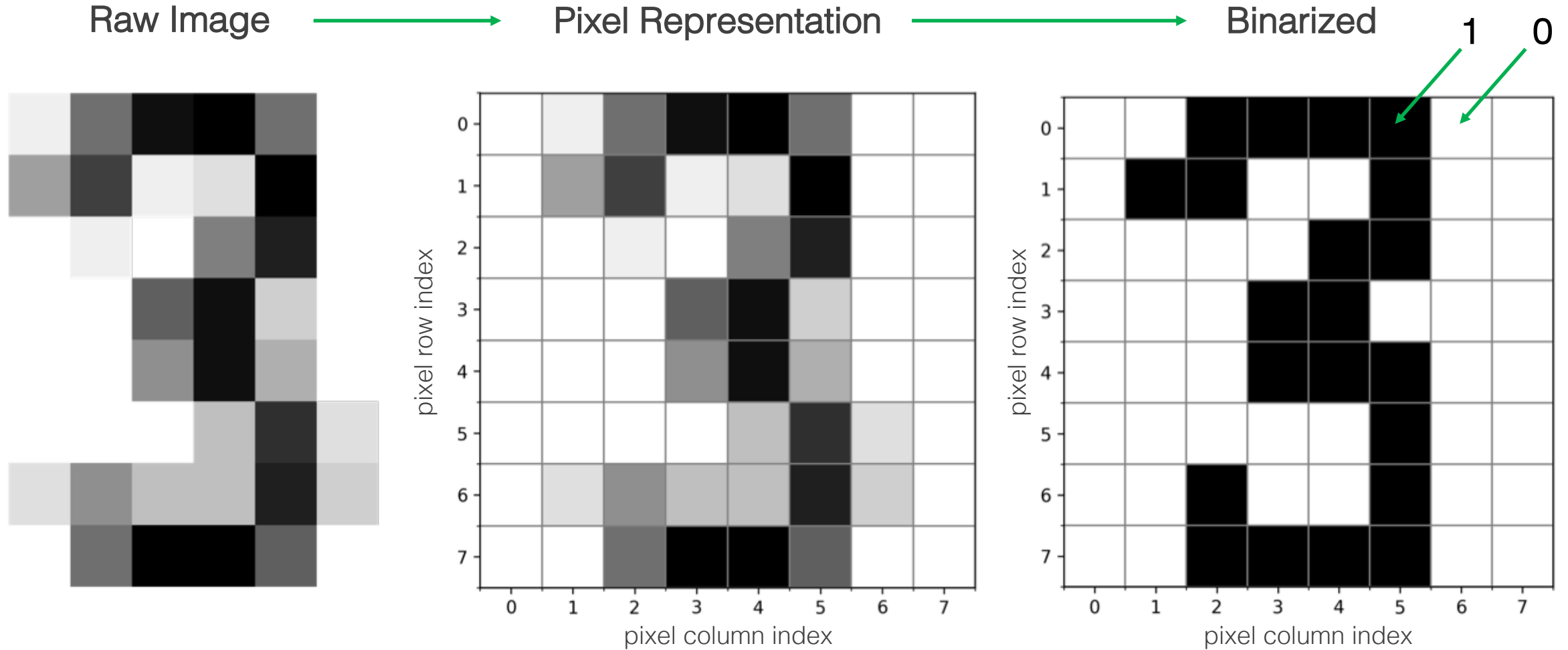
We search for  
the model that  
best fits our  
data



# Preprocessing

prepare your data for analysis

# How do we make a prediction rule?



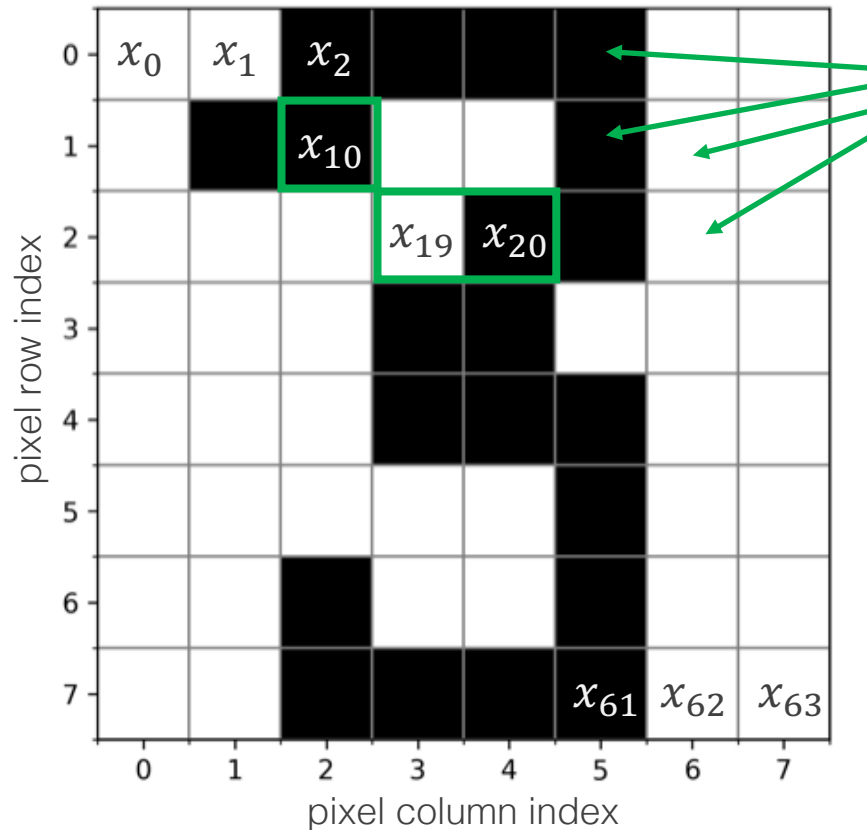


# Choose your model

Which prediction rules will you consider?

# How do we make a prediction rule?

$\mathbf{x}_i$ : one sample (observation) of the data



Each pixel represents an element of  $\mathbf{x}_i = [x_0, x_1, \dots, x_{63}]$

**64-dimensional binary data**

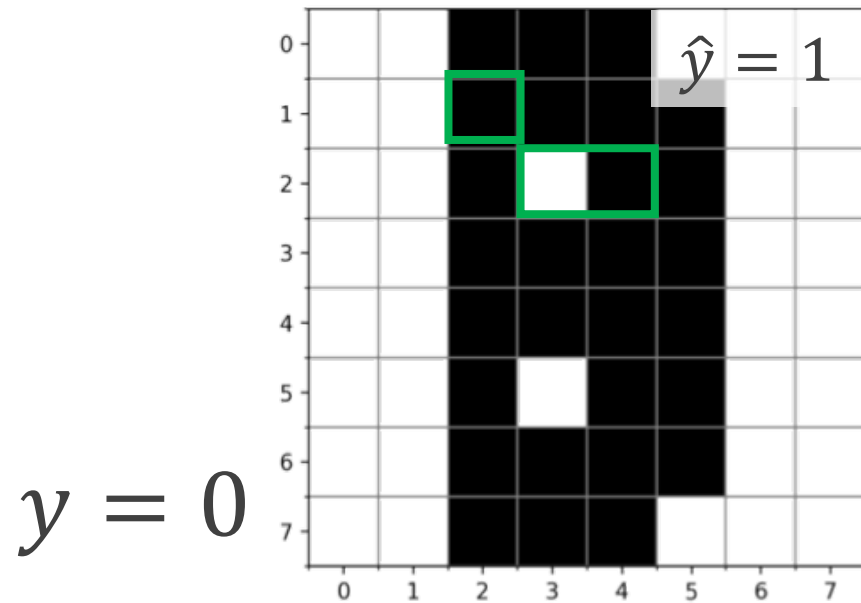
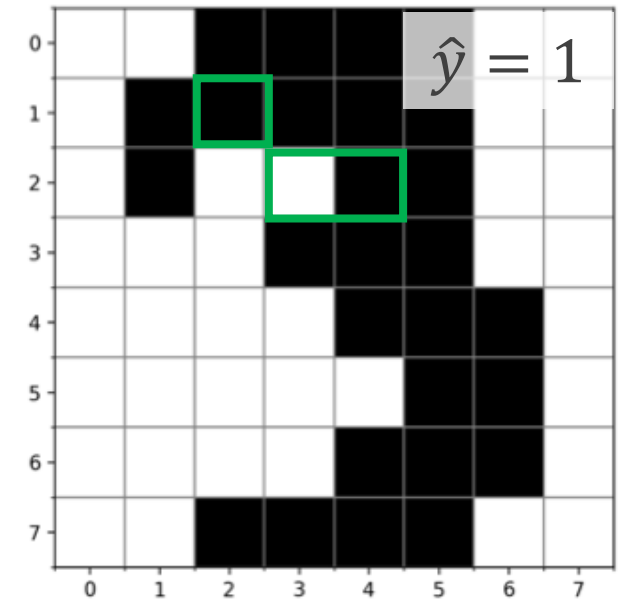
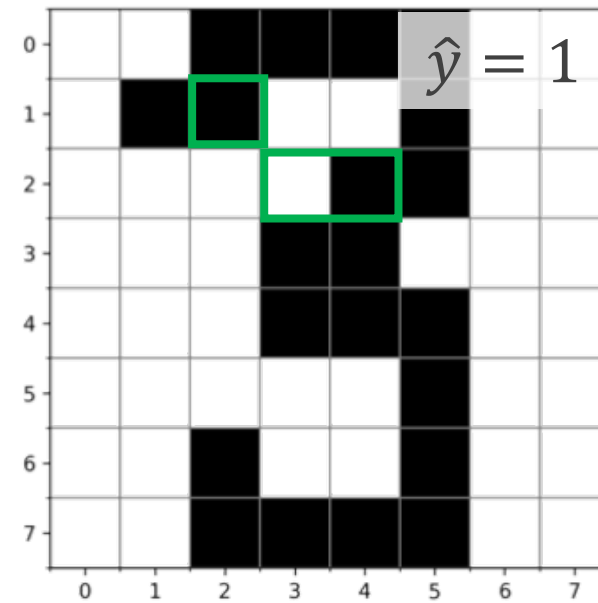
We want to estimate a function that takes these pixel values and decides whether this is a '3' or not

Example decision rule / prediction algorithm:

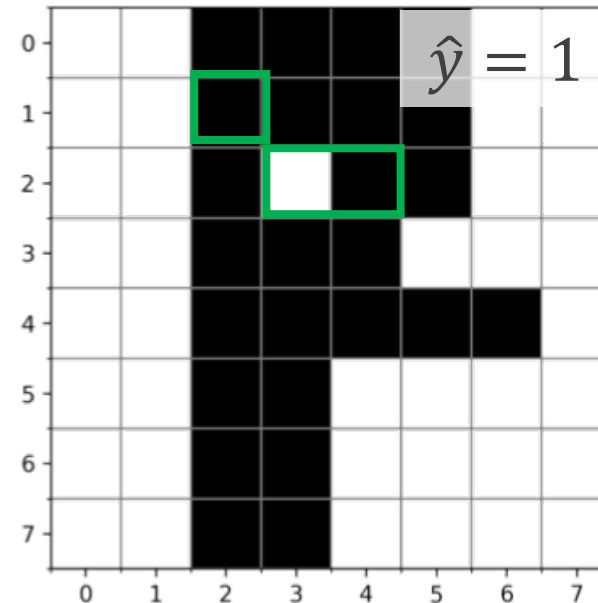
$$f(\mathbf{x}) = \begin{cases} 1 & \mathbf{x}_{10}\mathbf{x}_{20}(1 - \mathbf{x}_{19}) > 0 \\ 0 & \text{else} \end{cases}$$

# Does the rule work on other examples?

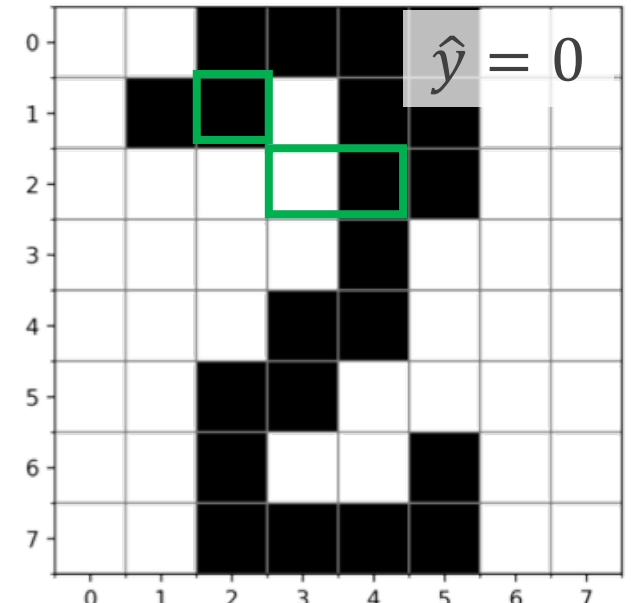
$$\hat{f}(x) = \begin{cases} 1 & x_{10}x_{20}(1 - x_{19}) > 0 \\ 0 & \text{else} \end{cases} \quad y = 1$$



✗

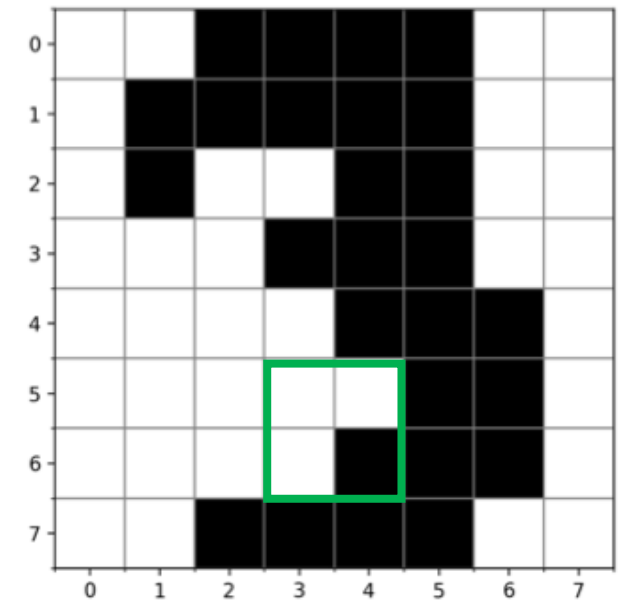
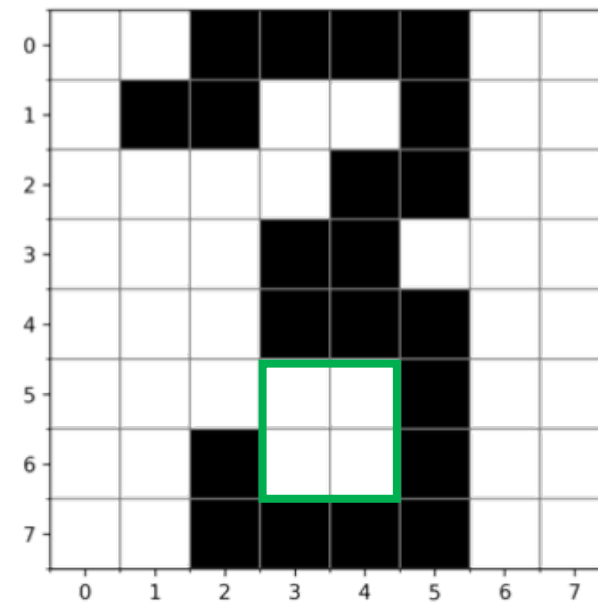


✗

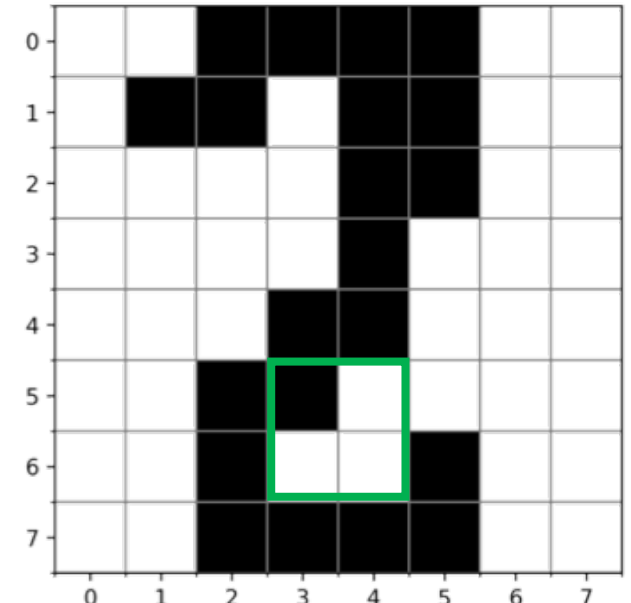
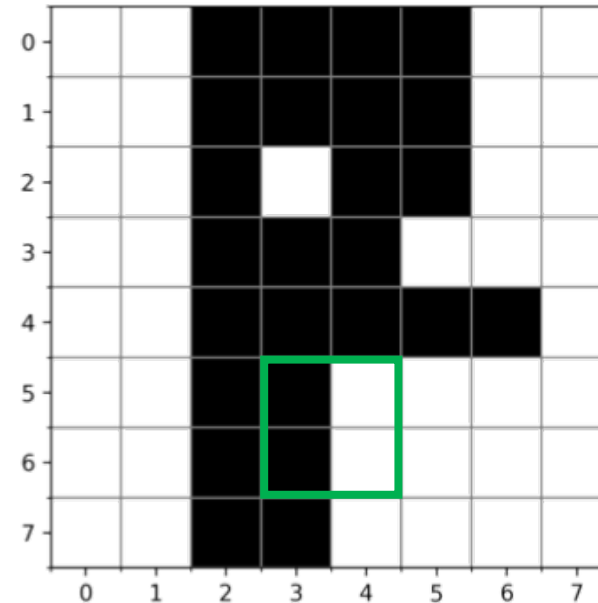
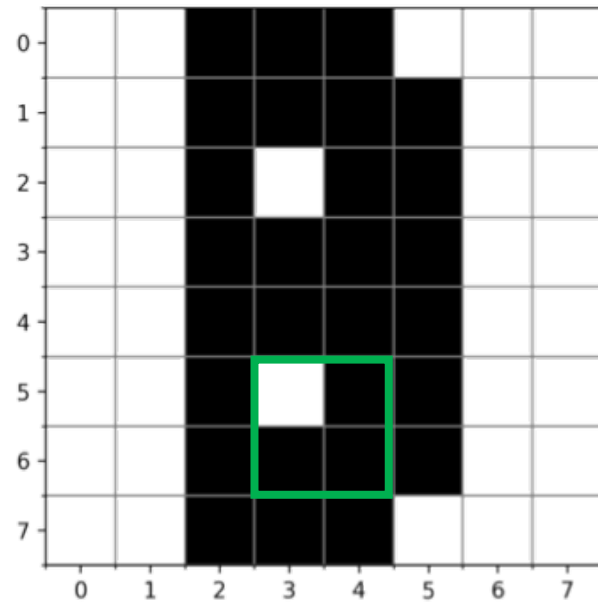


This is very hard to do manually, let's limit our search to a **subset** of the features...

$$y = 1$$



$$y = 0$$



# What decision rules work?

$y = 1$


$y = 0$


Option 1:

$$\hat{f}_1(\mathbf{x}) = \begin{cases} 1 & (1 - x_a)(1 - x_b)(1 - x_c) > 0 \\ 0 & \text{else} \end{cases}$$

Option 2:

$$\hat{f}_2(\mathbf{x}) = \begin{cases} 1 & (1 - x_a)(1 - x_b) > 0 \\ 0 & \text{else} \end{cases}$$

For simplicity, we refer to these  $x_i$  values as:

$x_a$	$x_b$
$x_c$	$x_d$

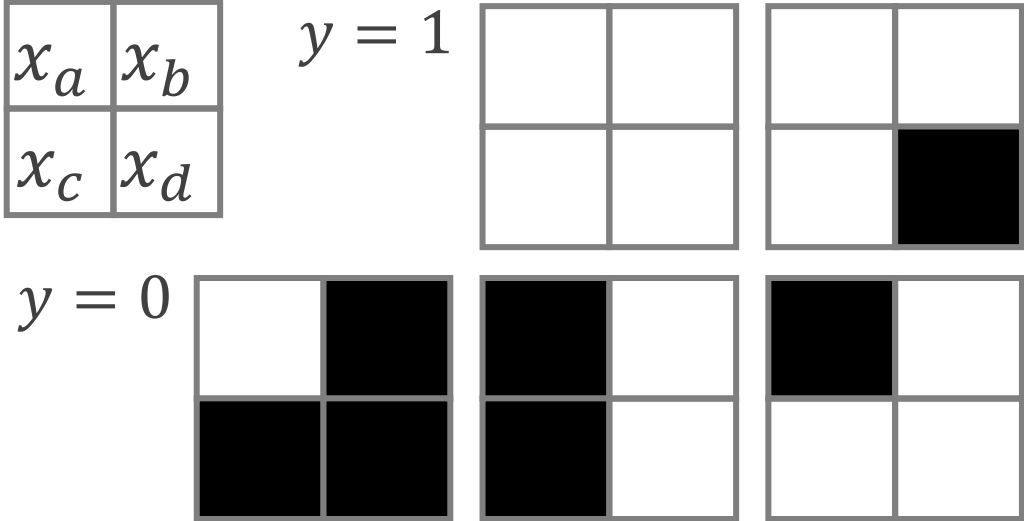
Option 3 (exclude our negative examples):

$$\hat{f}_3(\mathbf{x}) = \begin{cases} 1 & 1 - [(1 - x_a)x_bx_cx_d + x_a(1 - x_b)x_c(1 - x_d) + x_a(1 - x_b)(1 - x_c)(1 - x_d)] > 0 \\ 0 & \text{else} \end{cases}$$

# Learning: train your model

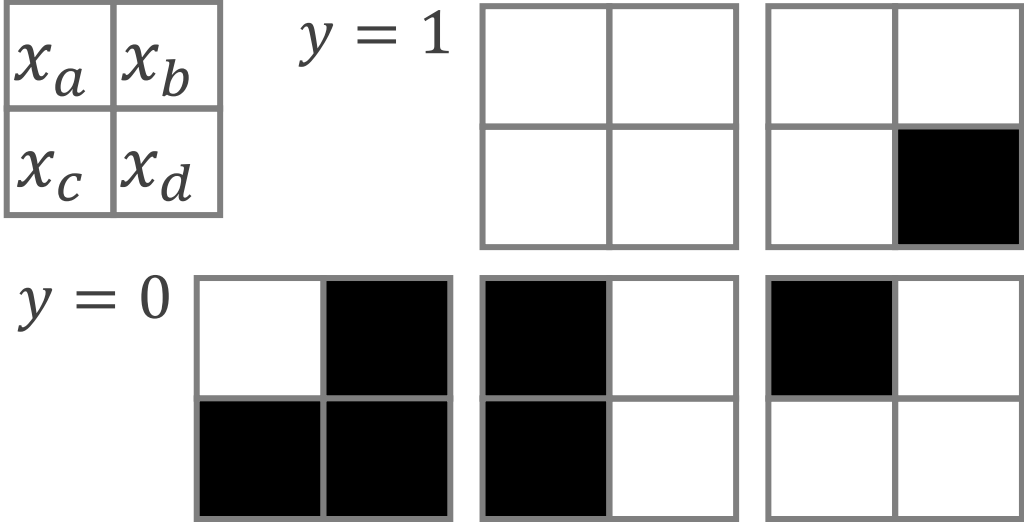
Select the best rule from the options available

$x_a$	$x_b$	$x_c$	$x_d$	$y$
0	0	0	0	<b>1</b>
0	0	0	1	<b>1</b>
0	0	1	0	
0	0	1	1	
0	1	0	0	
0	1	0	1	
0	1	1	0	
0	1	1	1	<b>0</b>
1	0	0	0	<b>0</b>
1	0	0	1	
1	0	1	0	<b>0</b>
1	0	1	1	
1	1	0	0	
1	1	0	1	
1	1	1	0	
1	1	1	1	



What are all the possible combinations?

$x_a$	$x_b$	$x_c$	$x_d$	$y$	$\hat{f}_1$	$\hat{f}_2$	$\hat{f}_3$
0	0	0	0	1	1	1	1
0	0	0	1	1	1	1	1
0	0	1	0		0	1	1
0	0	1	1		0	1	1
0	1	0	0		0	0	1
0	1	0	1		0	0	1
0	1	1	0		0	0	1
0	1	1	1	0	0	0	0
1	0	0	0	0	0	0	0
1	0	0	1		0	0	1
1	0	1	0	0	0	0	0
1	0	1	1		0	0	1
1	1	0	0		0	0	1
1	1	0	1		0	0	1
1	1	1	0		0	0	1
1	1	1	1		0	0	1



Option 1:

$$\hat{f}_1(\boldsymbol{x}) = 1 \text{ if } (1 - x_a)(1 - x_b)(1 - x_c) > 0$$

Option 2:

$$\hat{f}_2(\boldsymbol{x}) = 1 \text{ if } (1 - x_a)(1 - x_b) > 0$$

Option 3 (exclude our negative examples):

$$\hat{f}_3(\boldsymbol{x}) = 1 \text{ if } 1 - [(1 - x_a)x_bx_cx_d + x_a(1 - x_b)x_c(1 - x_d) + x_a(1 - x_b)(1 - x_c)(1 - x_d)] > 0$$



Possible observations =  $2^n = 16$

Features:  $n = 4$

Decision rules:  $N = 2^{2^n} = 65,536$

$x_a$	$x_b$	$x_c$	$x_d$	$y$	$\hat{h}_0$	$\hat{h}_1$	$\hat{h}_2$	...	$\hat{h}_{N-1}$	$\hat{h}_N$
0	0	0	0	1	0	1	0		1	1
0	0	0	1	1	0	0	1		1	1
0	0	1	0		0	0	0		1	1
0	0	1	1		0	0	0		1	1
0	1	0	0		0	0	0		1	1
0	1	0	1		0	0	0		1	1
0	1	1	0		0	0	0		1	1
0	1	1	1	0	0	0	0		1	1
1	0	0	0	0	0	0	0		1	1
1	0	0	1		0	0	0		1	1
1	0	1	0	0	0	0	0		1	1
1	0	1	1		0	0	0		1	1
1	1	0	0		0	0	0		1	1
1	1	0	1		0	0	0		1	1
1	1	1	0		0	0	0		1	1
1	1	1	1		0	0	0		0	1

We have  $n = 4$  features:

$x_a, x_b, x_c, x_d$

They produce  $2^n = 16$   
**possible observations** (i.e.  
unique feature combinations)...

...yielding  $N = 2^{2^4} = 65,536$   
**unique functions** to predict  $y$ !

There are 5 labeled  
observations, so  $16 - 5 = 11$   
unlabeled observations. So  
there are  $2^{11} = 2,048$  **unique  
functions that all correctly  
predict  $y$**  for our 5 observations

# We cannot try **all possible decision rules** for features...

Number of binary features

Unique Decision Rules

1	$2^{2^1} =$	<b>4</b>
2	$2^{2^2} =$	<b>16</b>
3	$2^{2^3} =$	<b>256</b>
4	$2^{2^4} =$	<b>65,536</b>
5	$2^{2^5} =$	<b>4,294,967,296</b>
6	$2^{2^6} =$	<b>18,446,744,073,709,551,616</b>
7	$2^{2^7} =$	<b>340,282,366,920,938,463,463,374,607,431,768,211,456</b>

For our digits data, we have **64 features**...

← 1 year of processing on the world's most powerful computer →

For continuous features, the number of unique rules is **infinite**

# Key Challenges

Manually selecting rules is typically impractical

We **cannot test all possible decision rules**

How do we choose the best rule given the data?

# Possible Solutions

Supervised learning searches a **subset** of all possible decision rules to “learn” (choose) a rule

Pick the one that **generalizes** best – often **simpler** (regularized) rules. This may mean allowing for some mistakes (tradeoff).

# Evaluate performance

See how well your learned model works on new data (i.e. generalizes)

# Let's test some rules

We divide the data we have into **training** and **validation** sets so we **never** train and test on the same data

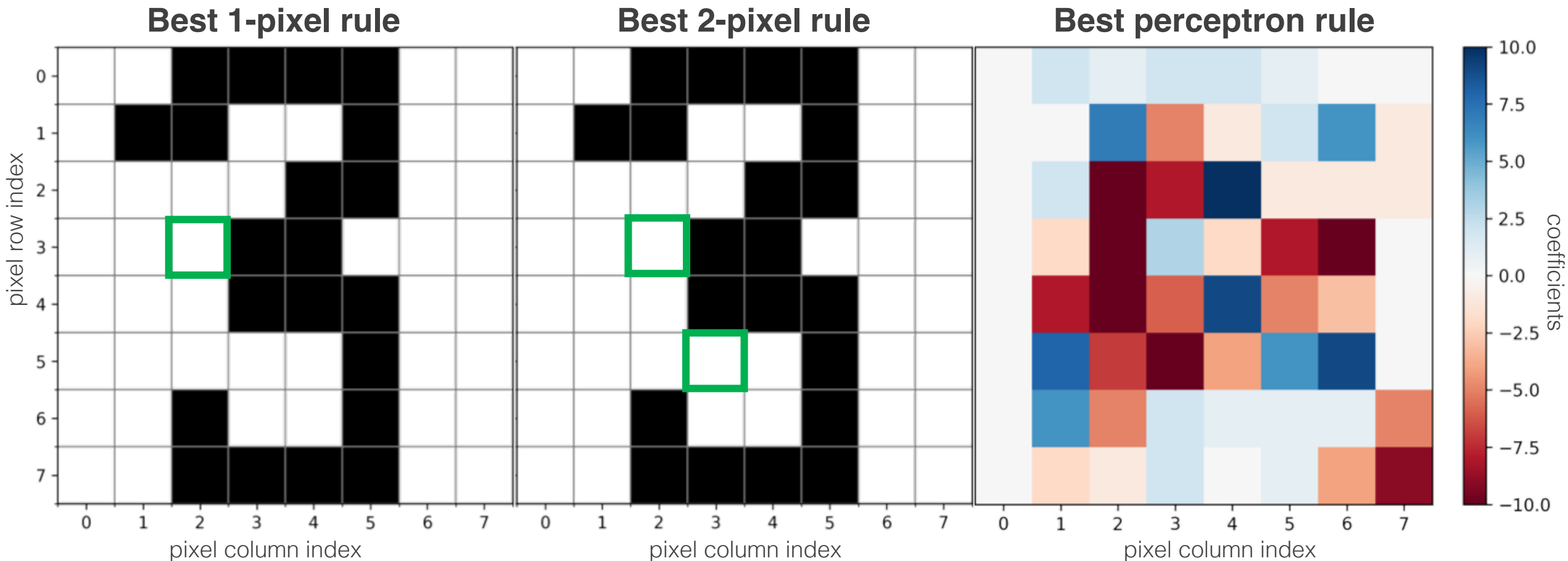
We limit our search for rules to three **hypothesis sets** (possible rules to search) using different **learning algorithms** (search for best rule)

Hypothesis sets (of rules)	Learning algorithm	Max # of features in rule
1. The best <b>one-pixel</b> rule	exhaustive search	1
2. The best <b>two-pixel</b> rule	exhaustive search	2
3. A <b>perceptron</b> rule (linear classifier)	stochastic gradient descent	64




**models**

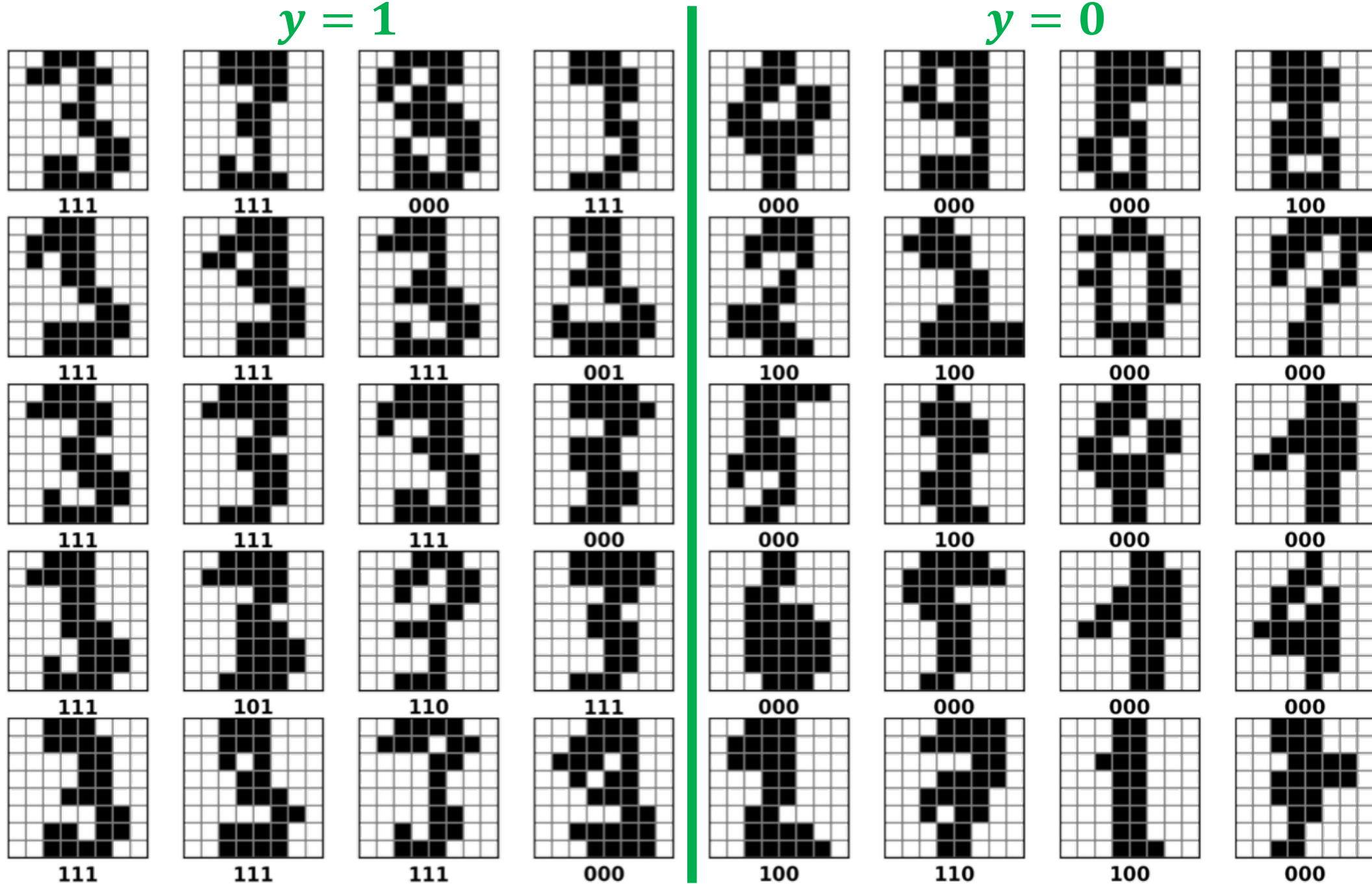
**optimization strategy**

# Our results:



# Example Results

 1-pixel rule  
 2-pixel rule  
 perceptron  
**010 Key**



# Summarizing supervised learning

Let's review and bring it all together



# Components of supervised learning

**Input**

$x$

**Output**

$y$

**Training Data**

$(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$

**Target function**

$f(x) \rightarrow y$

This is unknown, but the best you could ever do

**Hypothesis set**

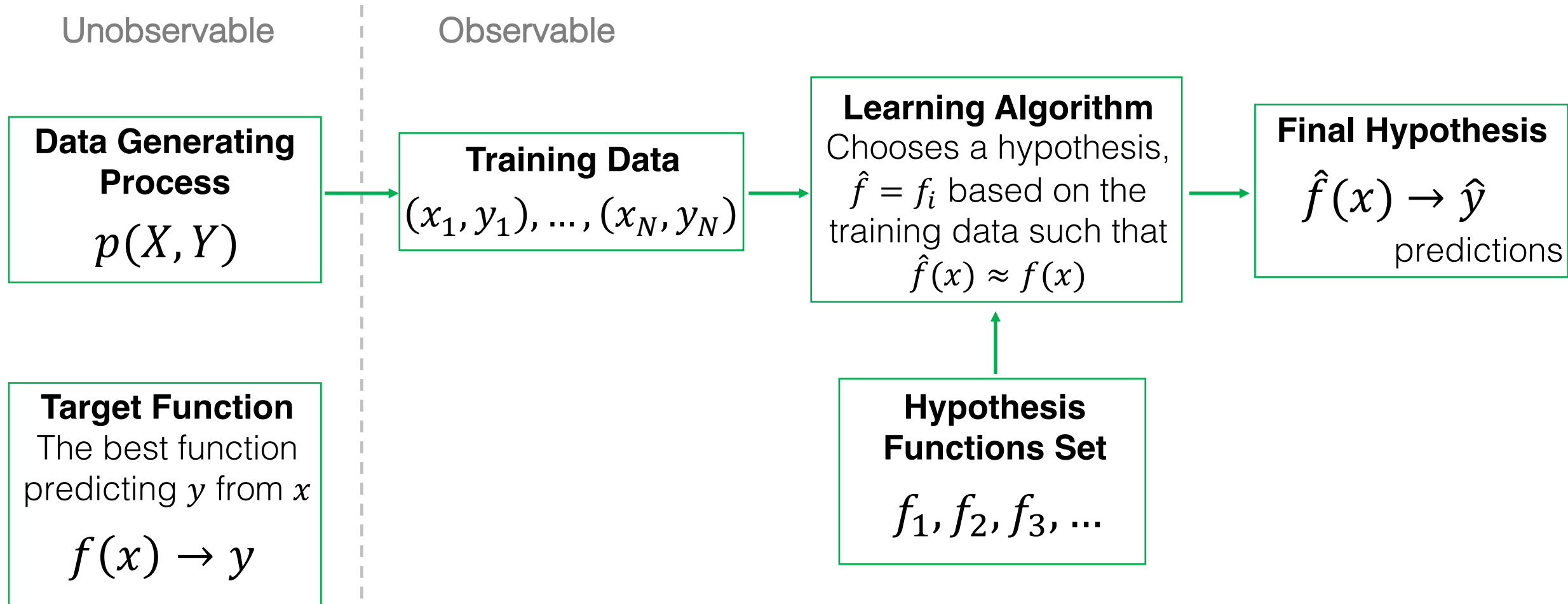
$f_i(x) \rightarrow \hat{y}$

Functions to consider in trying to approximate  $f(x)$

**Learning algorithm**

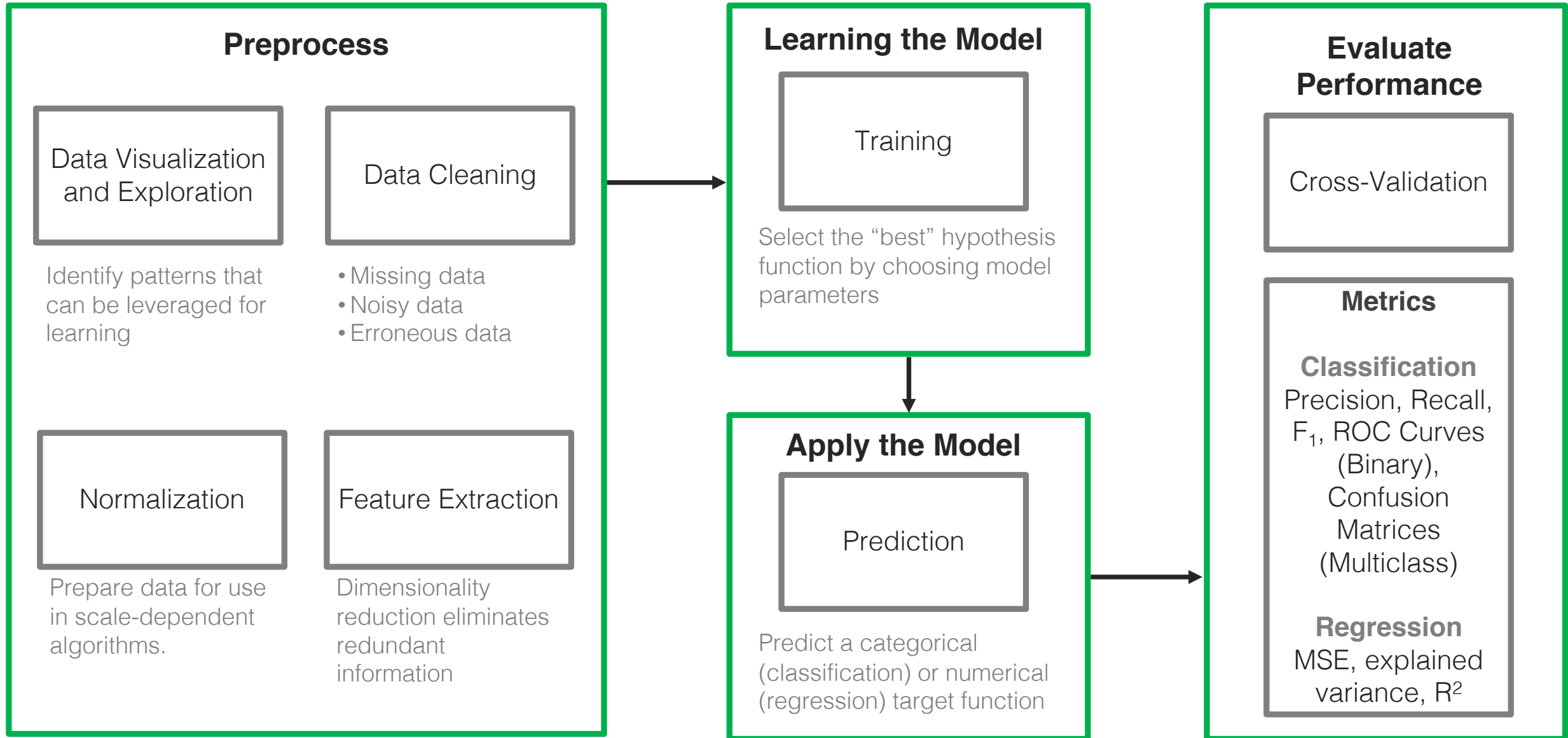
Optimization technique that searches the hypothesis set for the function  $f_i$  that best approximates  $f$  (typically by choosing parameters in a model)

# Supervised Learning



- Need to select the hypothesis functions (models to train)
- Need to select the learning algorithm (for fitting the models to the data)

# Supervised learning in practice



# References

## Further reading:

Abu-Mostafa, Yaser S., Malik Magdon-Ismail, and Hsuan-Tien Lin. Learning from data. Vol. 4. New York, NY, USA:: AMLBook, 2012.

Friedman, Jerome, Trevor Hastie, and Robert Tibshirani. The elements of statistical learning. Vol. 1. New York: Springer series in statistics, 2001.

Moore, Cristopher, and Stephan Mertens. The nature of computation. OUP Oxford, 2011.

## Videos/presentations that inspired this lecture:

[Learning from Data](#), Yaser Abu-Mostafa, Caltech

[Learning to See](#), Welch Labs