

Introduction to Machine Learning (and Notation)

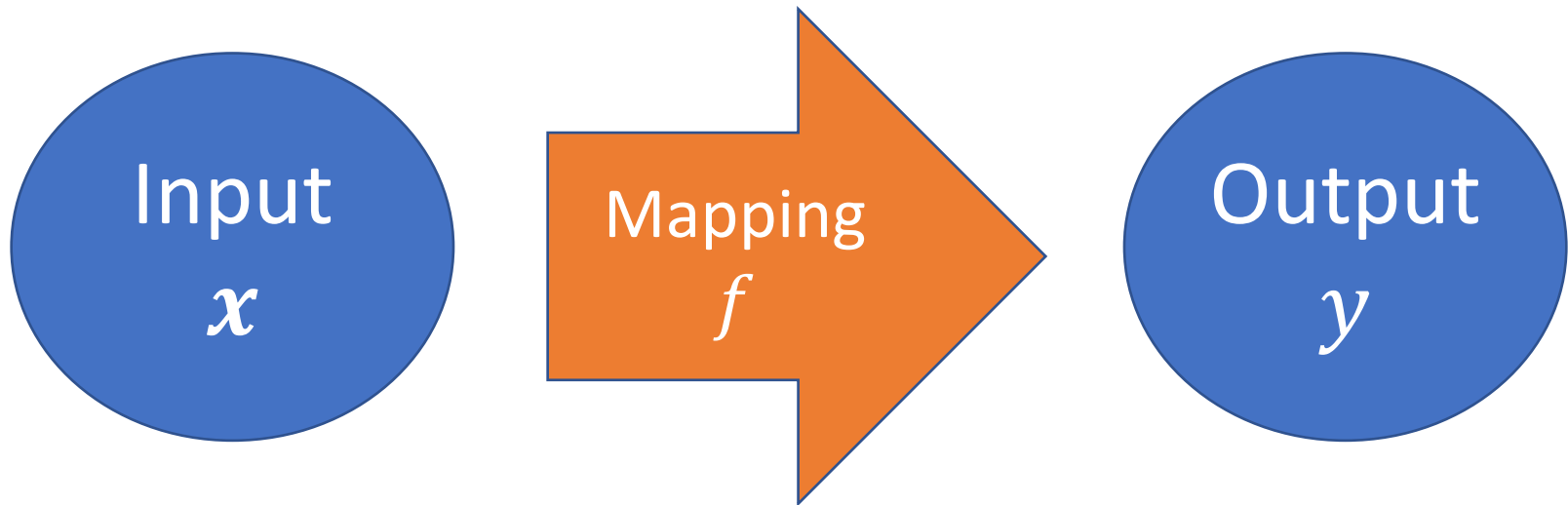
David I. Inouye

Friday, September 4, 2020

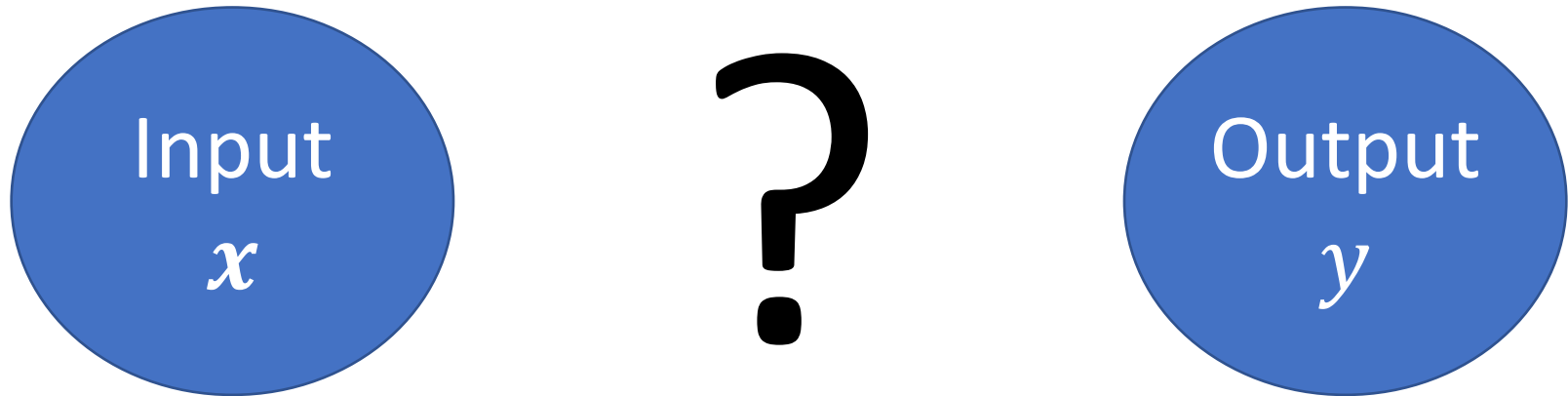
Outline

- ▶ Supervised learning
 - ▶ Regression
 - ▶ Classification
- ▶ Unsupervised learning
 - ▶ Dimensionality reduction (PCA)
 - ▶ Clustering
 - ▶ Generative models
- ▶ Other key concepts
 - ▶ Generalization
 - ▶ Curse of dimensionality
 - ▶ No free lunch theorem

The goal of supervised learning is to estimate a **mapping (or function)** between input and output



The goal of supervised learning is to estimate a **mapping (or function)** between input and output *given only input-output examples*



The set of input-output pairs is called a training set, denoted by $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$

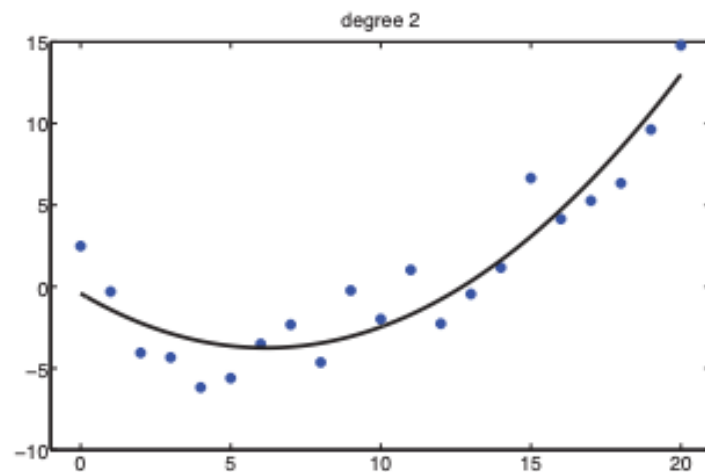
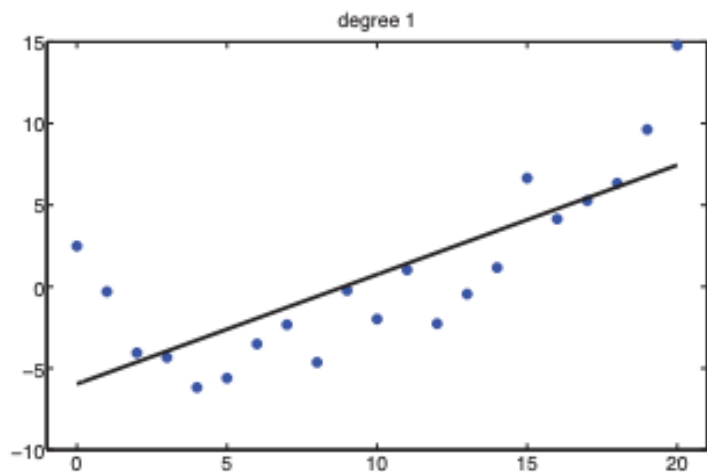
▶ Input \mathbf{x}_i

- ▶ Called features (ML), attributes, or covariates (Stats). Sometimes just variables.
- ▶ Can be numeric, categorical, discrete, or nominal.
- ▶ Examples
 - ▶ [height, weight, age, gender]
 - ▶ $[x_1, x_2, \dots, x_d]$ – A d -dimensional vector of numbers
 - ▶ Image
 - ▶ Email message

▶ Output y_i

- ▶ Called output, response, or target (or label)
- ▶ Real-valued/numeric output: e.g., $y_i \in \mathcal{R}$
- ▶ Categorical, discrete, or nominal output: y_i from *finite* set, i.e., $y_i \in \{1, 2, \dots, c\}$

If the output y_i is numeric,
then the problem is known as regression



NOTE: Input x does not have to be numeric. Only the output y must be numeric.

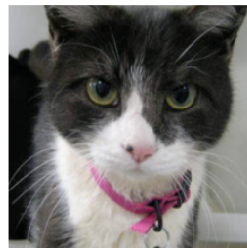
- ▶ Given height x_i , predict age y_i
- ▶ Predict GPA given SAT score
- ▶ Predict SAT score given GPA
- ▶ Predict GRE given SAT and GPA

If output is categorical,
then the problem is known as classification

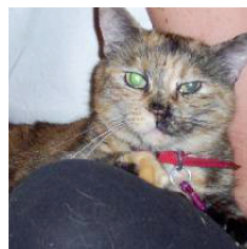
▶ Given height x ,
predict “male” ($y = 0$)
or “female” ($y = 1$)

▶ Given salary x_1 and
mortgage payment x_2 ,
predict defaulting on
loan (“yes” or “no”)

predicted: cat



predicted: cat



predicted: dog



predicted: cat



predicted: cat



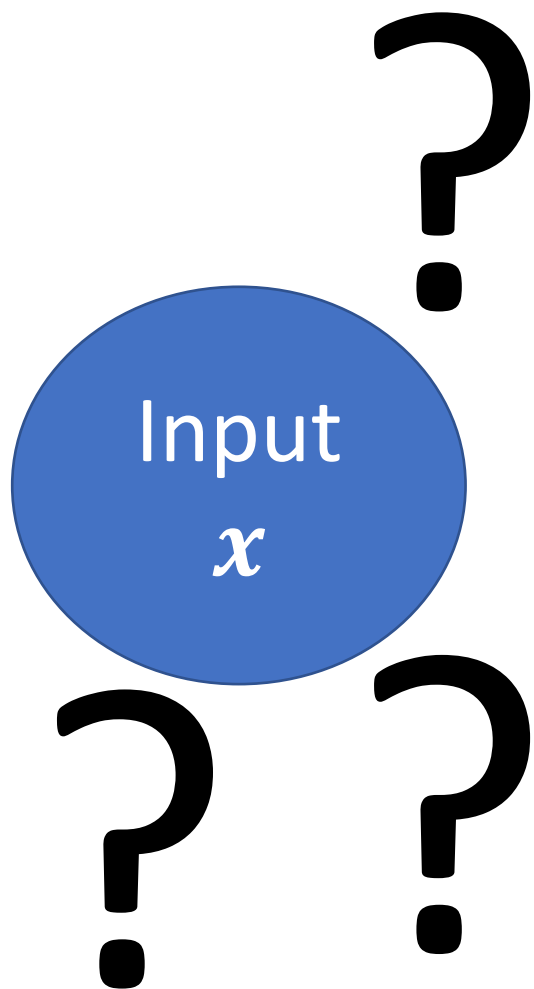
predicted: dog



Side note: Encoding / representing a categorical variable can be done in many ways

- ▶ Suppose the categorical variable is “yes” and “no”
 - ▶ Canonical ways: “no” \rightarrow 0 and “yes \rightarrow 1
 - ▶ What are other possible encodings?
- ▶ What if there are more than two categories such as cats, dogs, fish and snakes?
- ▶ What is good and bad about using $\{1,2,3,4\}$ for above example of animals?
- ▶ One-hot encoding is another common way

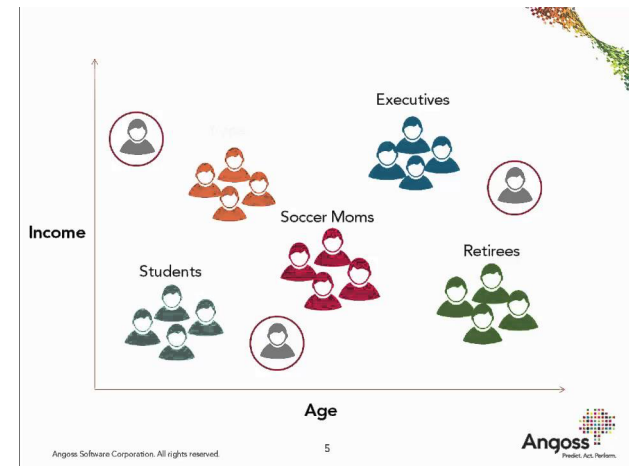
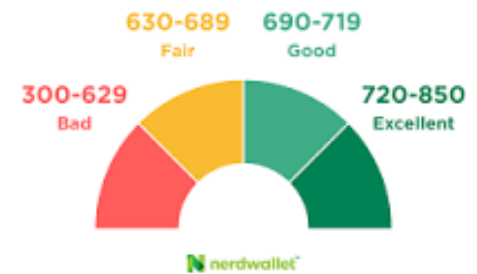
The goal of unsupervised learning is to model or understand the input data directly



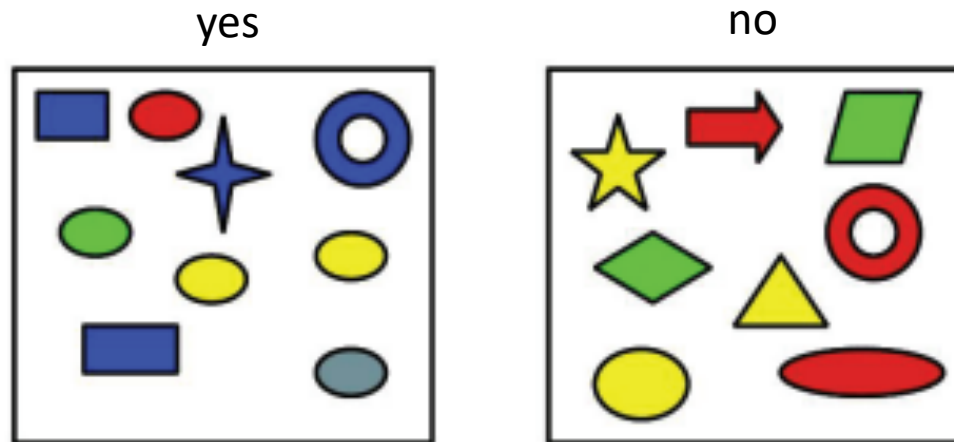
- ▶ Dimensionality reduction
- ▶ Clustering
- ▶ Generative models
“What I cannot create I do not understand”
– Richard Feynman

In unsupervised learning, the training set is only a set of input values $\mathcal{D} = \{\mathbf{x}_i\}_{i=1}^n$

- ▶ [Dimensionality reduction] Estimate a single number that summarizes all variables of wealth (e.g. credit score)
- ▶ [Clustering] Estimate natural groups of customers
- ▶ [Generative Models] Estimate the distribution of normal transactions to detect fraud (anomalies)



Given this dataset, should we use supervised or unsupervised learning?

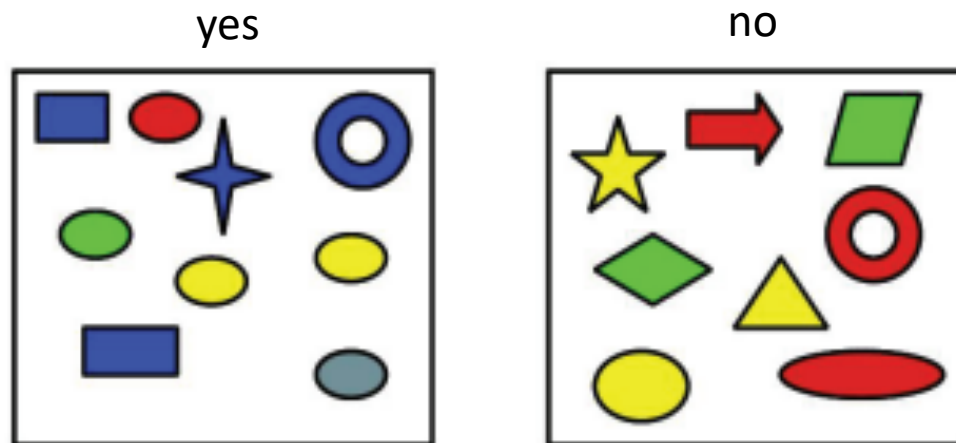


d features/attributes/covariates

n samples/
observations/
examples

Color	Shape	Size (cm)	Is it good?
Blue	Square	10	yes
Red	Ellipse	2.4	yes
Red	Ellipse	20.7	no

The dataset cannot determine the task, rather **the context** determines the task

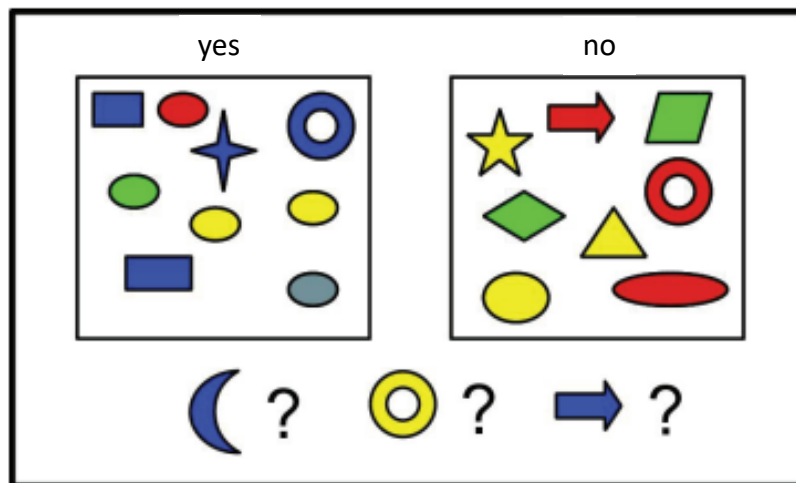


d features/attributes/covariates

n samples/
observations/
examples

Color	Shape	Size (cm)	Is it good?
Blue	Square	10	yes
Red	Ellipse	2.4	yes
Red	Ellipse	20.7	no

Generalization *beyond* the training set is the main goal of learning

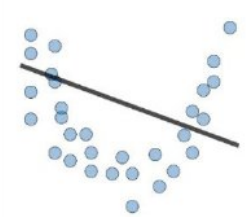

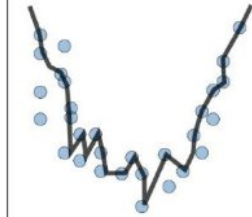
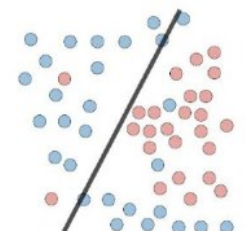
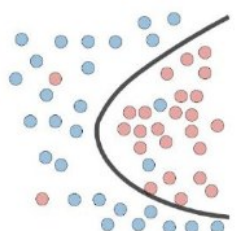
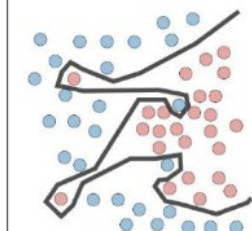
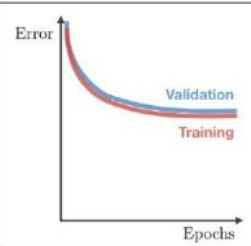
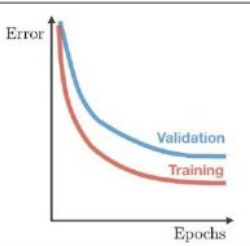
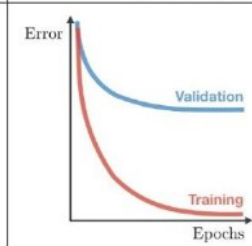


d features/attributes/covariates

		Color	Shape	Size (cm)	Is it good?		
n samples/ observations/ examples	x_1	Blue	Square	10	yes	y_1	
	x_2	Red	Ellipse	2.4	yes	y_2	
		Red	Ellipse	20.7	no		

Example from Machine Learning: A Probabilistic Perspective, Ch. 1, Kevin P. Murphy, 2012.

Generalization *beyond* the training set is the main goal of learning

	Underfitting	Just right	Overfitting
Symptoms	<ul style="list-style-type: none"> - High training error - Training error close to test error - High bias 	<ul style="list-style-type: none"> - Training error slightly lower than test error 	<ul style="list-style-type: none"> - Low training error - Training error much lower than test error - High variance
Regression			
Classification			
Deep learning			
Remedies	<ul style="list-style-type: none"> - Complexify model - Add more features - Train longer 		<ul style="list-style-type: none"> - Regularize - Get more data

Original source for figure unknown.

The curse of dimensionality is *unintuitive*

Example: Most space is in the “corners”

- ▶ Ratio between unit hypersphere to unit hypercube

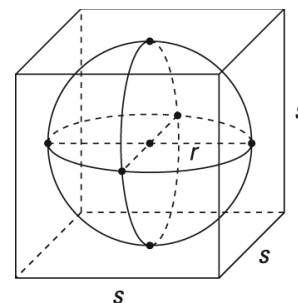
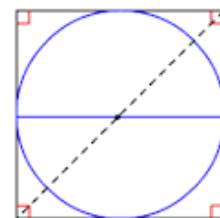
- ▶ 1D : $2/2 = 1$

- ▶ 2D : $\frac{\pi}{4} = 0.7854$

- ▶ 3D : $\frac{3\sqrt{3}\pi}{8} = 0.5238$

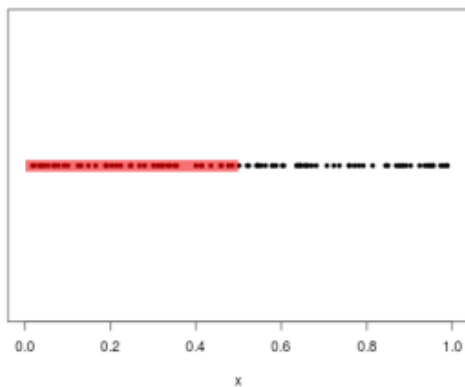
- ▶ d-dimensions: $V_d(r) = \frac{\pi^{\frac{n}{2}}}{\Gamma(\frac{n}{2}+1)} r^n$

- ▶ Thus, for 10-D: $2.55/2^{10} = 2.55/1024 = 0.00249$

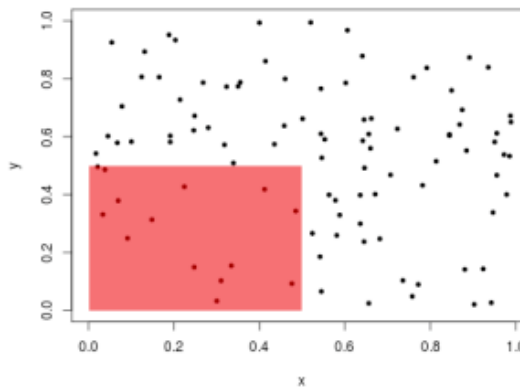


The curse of dimensionality is *unintuitive* *The number of points in 1/2 cube is very small*

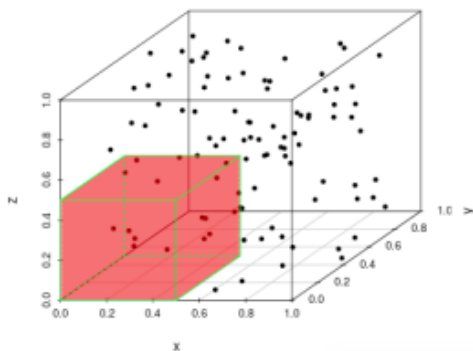
1-D: 42% of data captured.



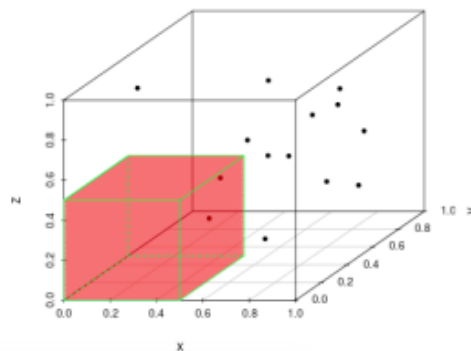
2-D: 14% of data captured.



3-D: 7% of data captured.



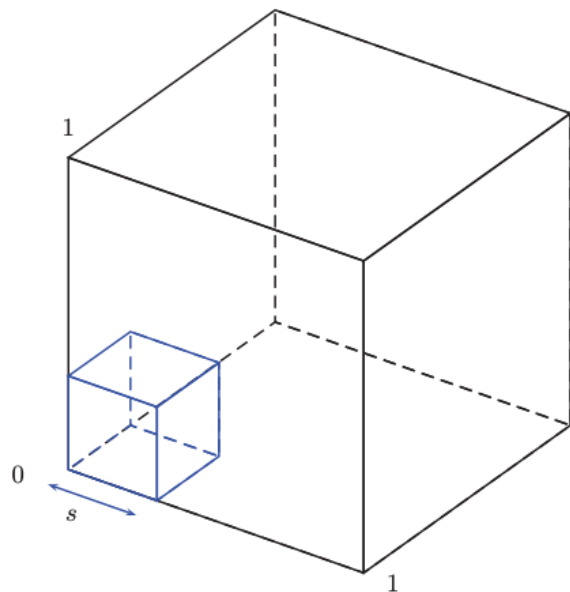
4-D: 3% of data captured.



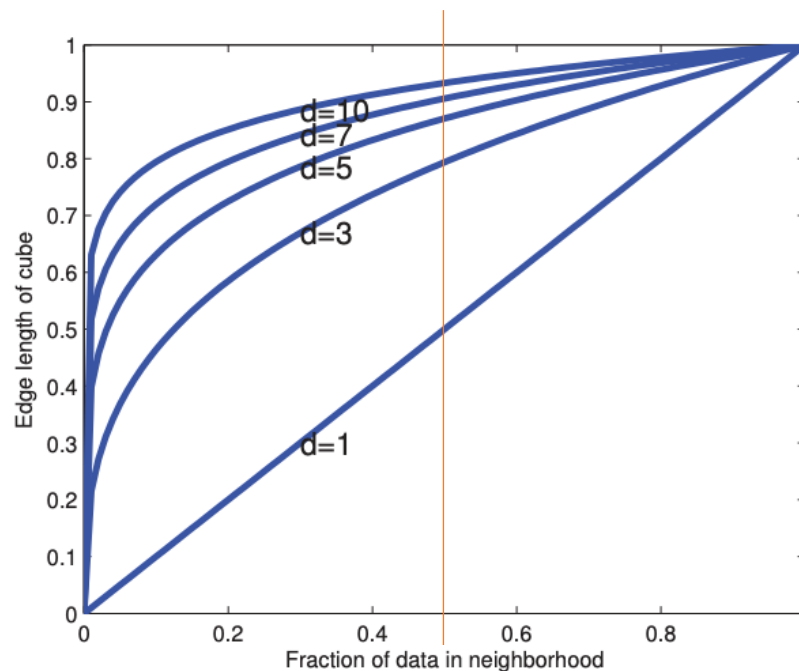
<https://eranraviv.com/curse-of-dimensionality/>

The curse of dimensionality is *unintuitive*

Example: Need edge length to be 0.9 to capture 1/2 data samples in 10 dimensions



(a)

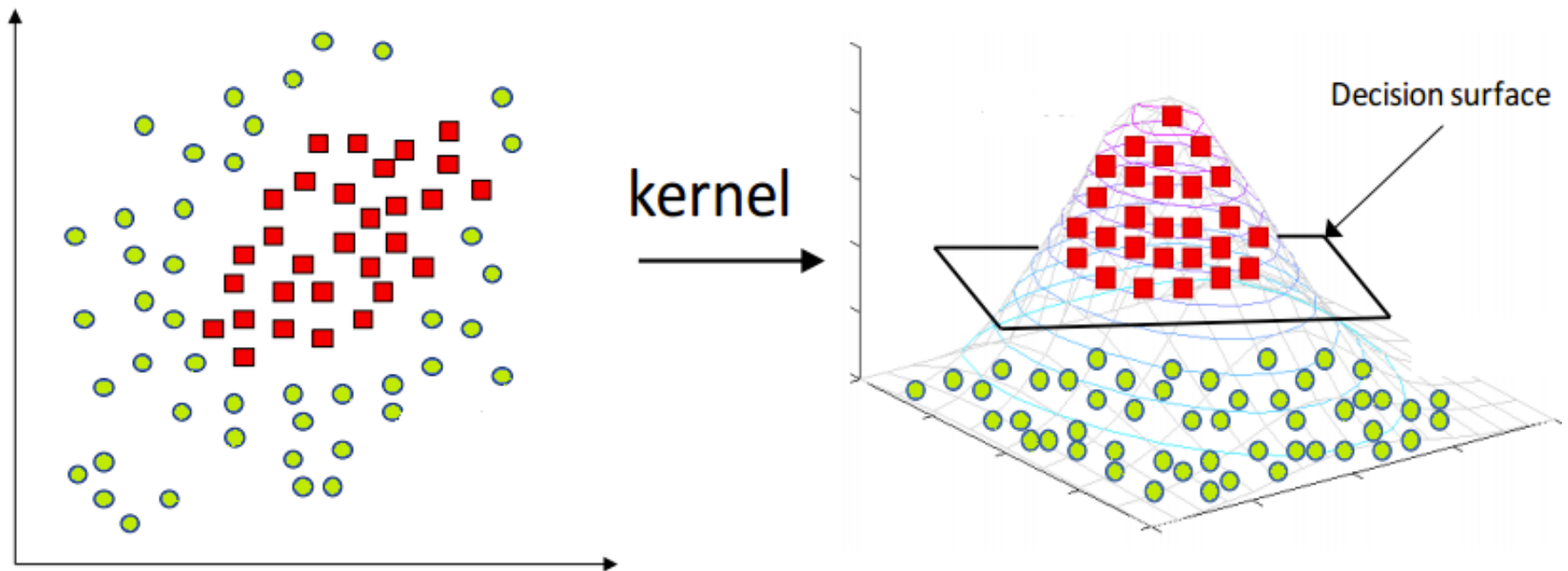


(b)

Figure 1.16 Illustration of the curse of dimensionality. (a) We embed a small cube of side s inside a larger unit cube. (b) We plot the edge length of a cube needed to cover a given volume of the unit cube as a function of the number of dimensions. Based on Figure 2.6 from (Hastie et al. 2009). Figure generated by `curseDimensionality`.

From *Machine Learning: A Probabilistic Perspective*, Kevin Murphy, 2012.

The “blessing” of dimensionality (more data generally doesn’t hurt if you can ignore)

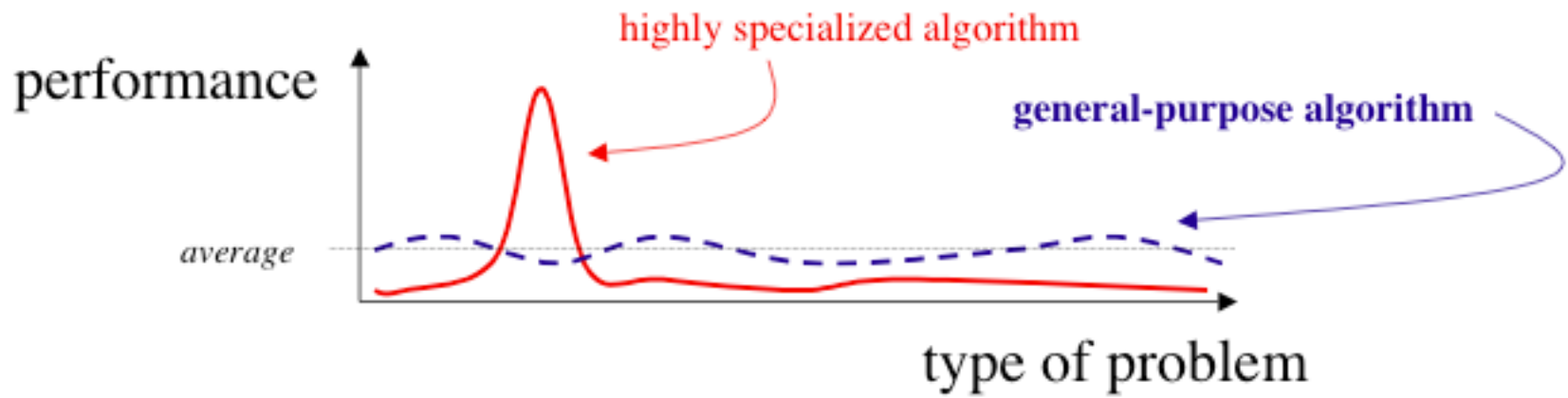


<https://www.hackerearth.com/blog/developers/simple-tutorial-svm-parameter-tuning-python-r/>

No Free Lunch Theorem

(“All models are wrong, but some models are useful.”*)

- ▶ All models are approximations
- ▶ All models make assumptions
- ▶ Assumptions are never perfect



* George Box (Box and Draper 1987, page 424).