
Machine Learning

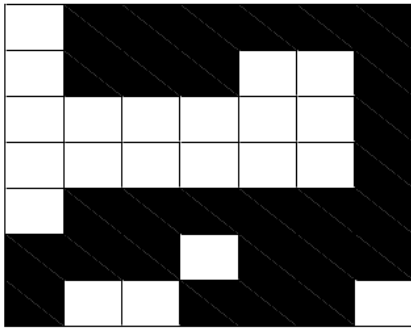
CSE 6363 (Fall 2016)

Lecture 2 Basic Machine Learning

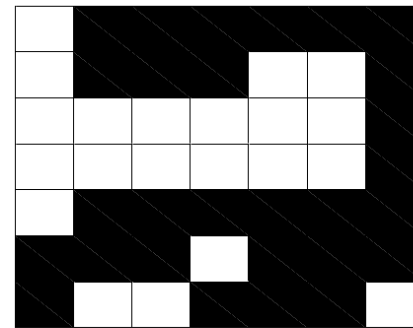
Heng Huang, Ph.D.

Department of Computer Science and Engineering

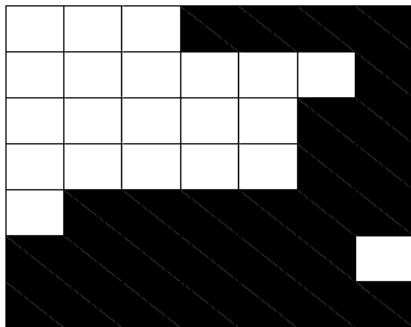
Learning, Biases, Representation



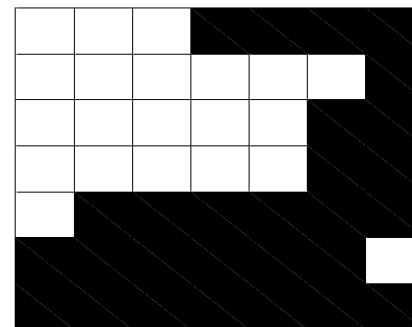
“yes”



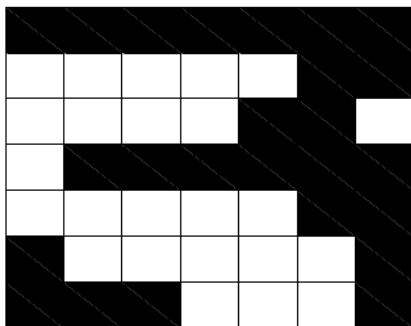
“yes”



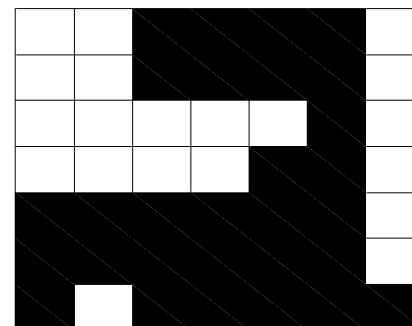
“yes”



“yes”



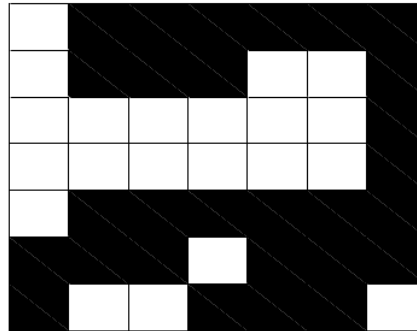
“no”



“no”

Representation

- There are many ways of presenting the same information



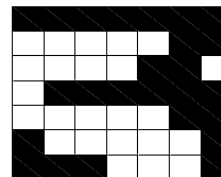
011111100111001000000010000000100111111011101111100
1110111110001

- The choice of representation may determine whether the learning task is very easy or very difficult

Hypothesis Class

- Representation: examples are binary vectors of length $d = 64$

$$\mathbf{x} = [111 \dots 0001]^T =$$



and labels $y \in \{-1, 1\}$ (“no”, “yes”)

- The mapping from examples to labels is a “linear classifier”

$$\hat{y} = \text{sign}(\theta \cdot \mathbf{x}) = \text{sign}(\theta_1 x_1 + \dots + \theta_d x_d)$$

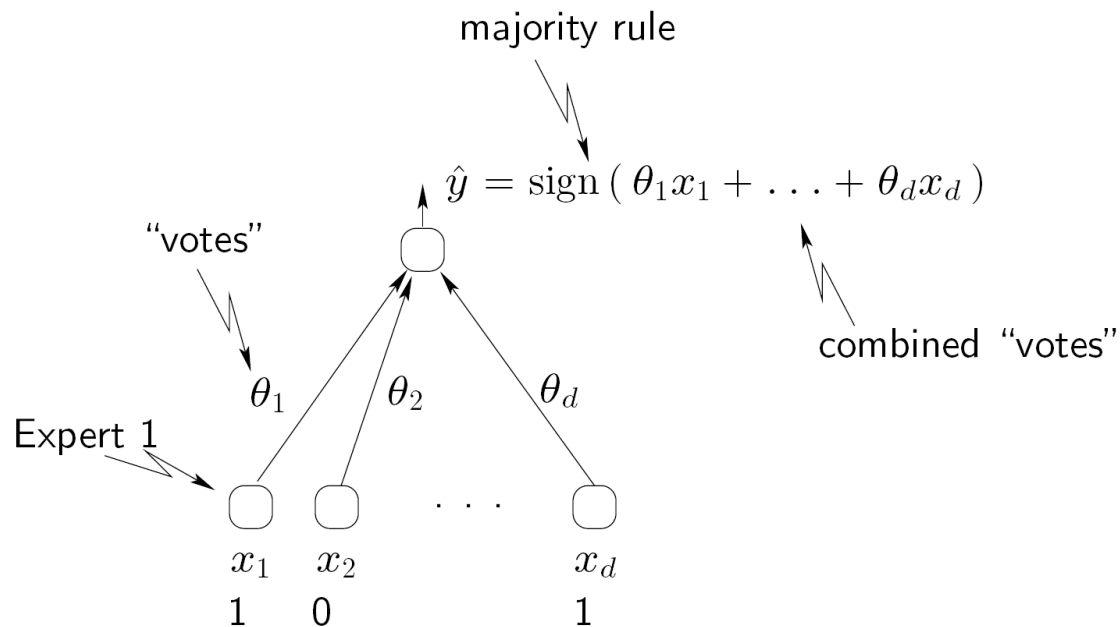
where θ is a vector of *parameters* we have to learn from examples.

Linear Classifier/Experts

- We can understand the simple linear classifier

$$\hat{y} = \text{sign}(\theta \cdot \mathbf{x}) = \text{sign}(\theta_1 x_1 + \dots + \theta_d x_d)$$

as a way of combining expert opinion (in this case simple binary features)



Estimation

\mathbf{x}	y
01111110011100100000000100000001001111110111011111001110111110001	+1
00011111000000011000001110000011001111110111111001111111100000011	+1
11111110000000110000011000111111000000111100000111110001101111111	-1
...	...

- How do we adjust the parameters θ based on the labeled examples?

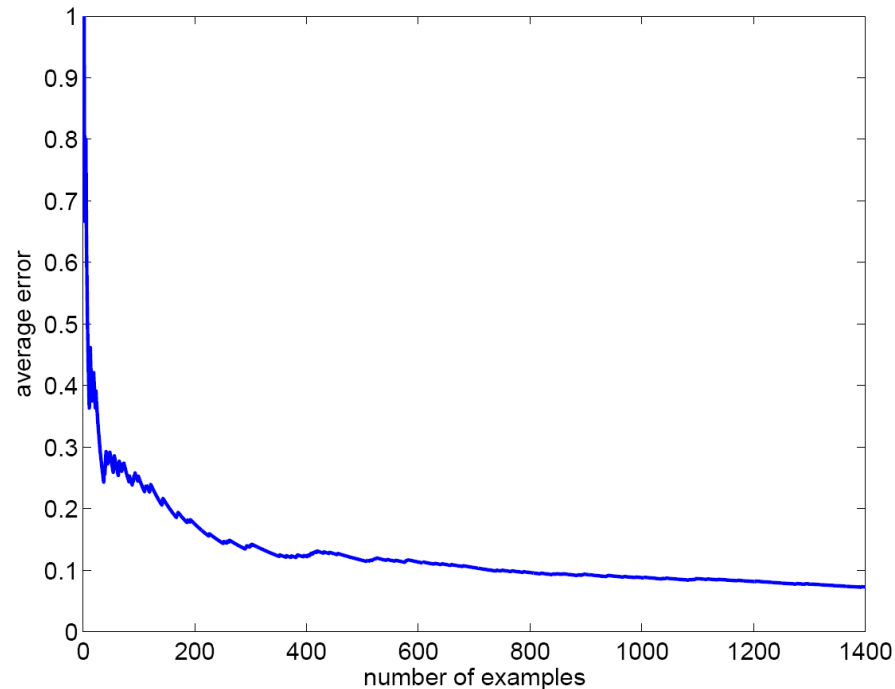
$$\hat{y} = \text{sign}(\theta \cdot \mathbf{x})$$

For example, we can simply refine/update the parameters whenever we make a mistake:

$$\theta_i \leftarrow \theta_i + y x_i, \quad i = 1, \dots, d \quad \text{if prediction was wrong}$$

Evaluation

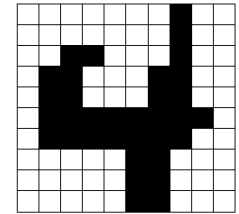
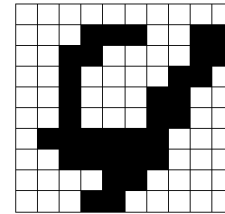
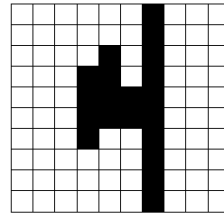
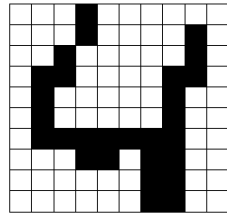
- Does the simple mistake driven algorithm work?



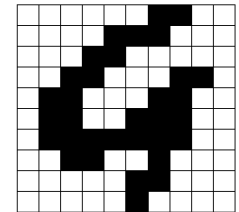
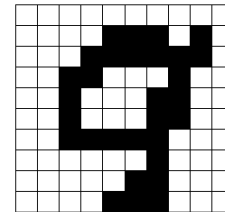
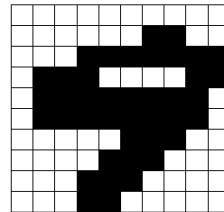
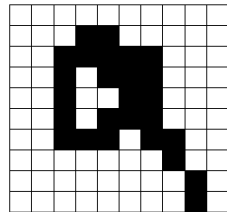
(average classification error as a function of the number of examples and labels seen so far)

Similar Problem

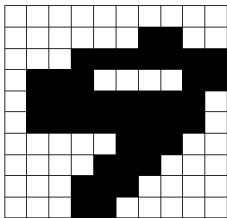
$y = +1$



$y = -1$



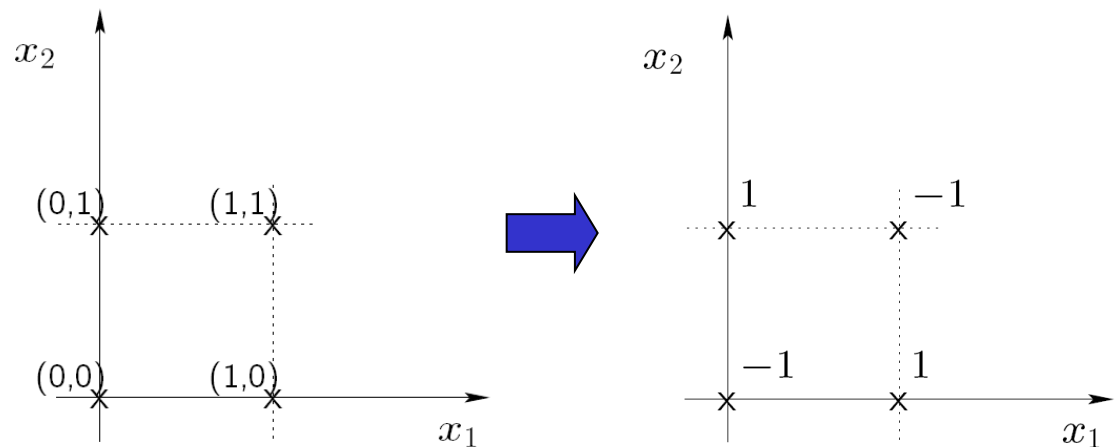
- Representation as a vector:



$$\Rightarrow [0000000000 \ 0000001100 \ 0001111111 \ \dots \ 0001100000]^T$$

Model Selection

- The simple linear classifier cannot solve all the problems (e.g., XOR)

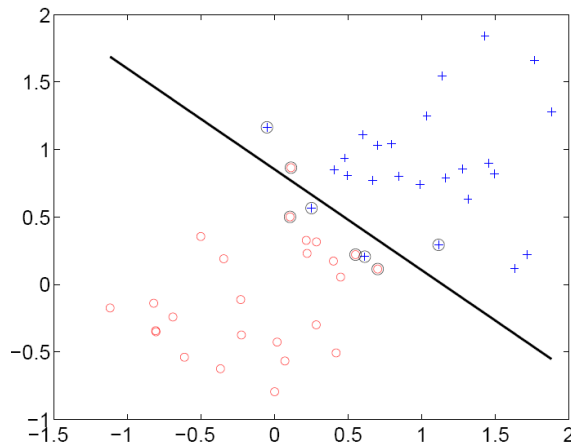


- Can we rethink the approach to do even better?

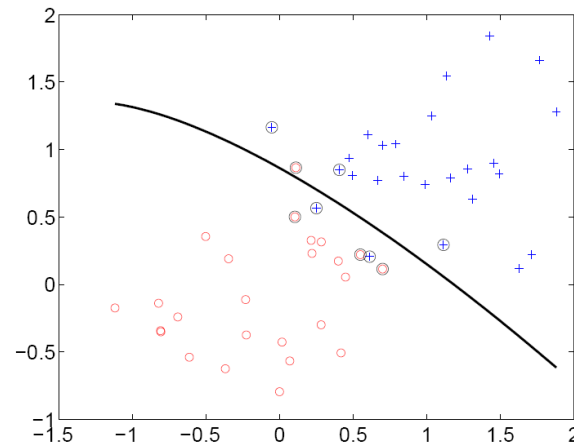
We can, for example, add “polynomial experts”

$$\hat{y} = \text{sign} (\theta_1 x_1 + \dots + \theta_d x_d + \theta_{12} x_1 x_2 + \dots)$$

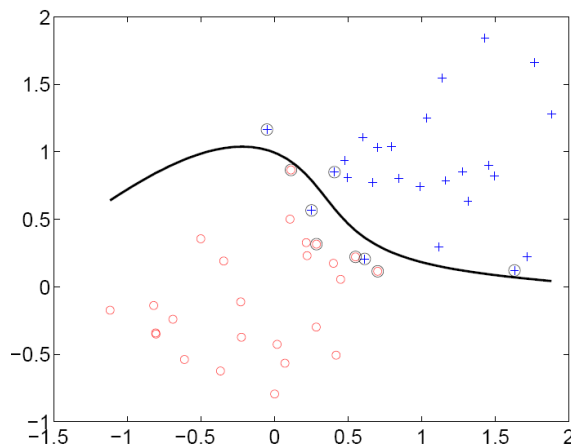
Model Selection (cont.)



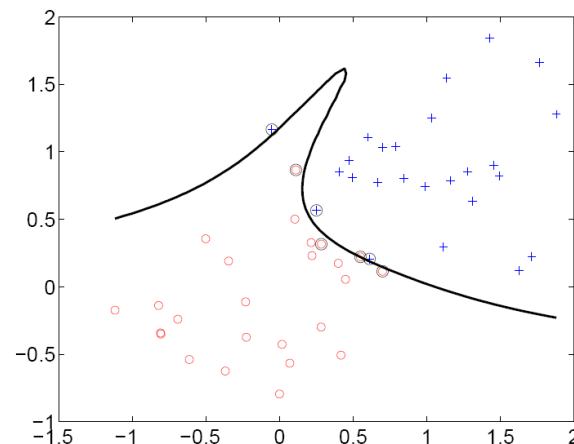
linear



2nd order polynomial



4th order polynomial

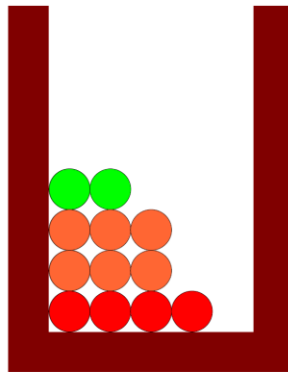


8th order polynomial

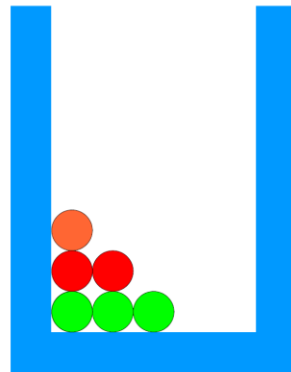
Probability Theory

- Boxes of fruit

● apple
● orange
● strawberry



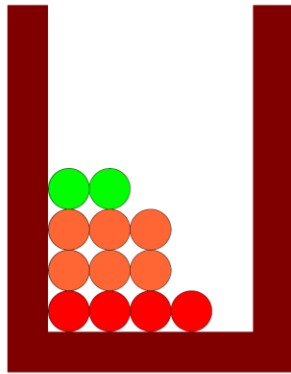
	apples	oranges	strawberries
red jar	2	6	4



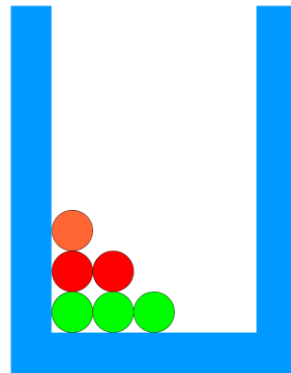
	apples	oranges	strawberries
blue jar	3	1	2

Probabilities of Fruit from a Given Jar

- apple
- orange
- strawberry



	apples	oranges	strawberries	
red jar	$\frac{2}{12}$ = 0.167	$\frac{6}{12}$ = 0.5	$\frac{4}{12}$ = 0.33	sum = 1.0



	apples	oranges	strawberries	
blue jar	$\frac{3}{6}$ = 0.5	$\frac{1}{6}$ = 0.167	$\frac{2}{6}$ = 0.33	sum = 1.0

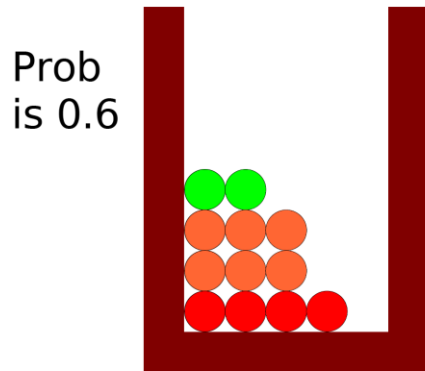
Choose Jar then Draw a Fruit

- apple
- orange
- strawberry

Say the probability of choosing a jar is

$$P(\text{Jar} = \text{red}) = 0.6$$

$$P(\text{Jar} = \text{blue}) = 0.4$$



The probability of choosing the red jar
and drawing an apple out of it is

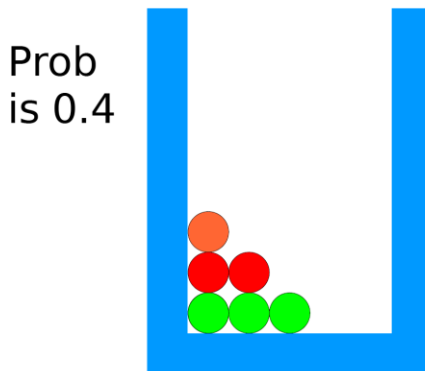
$$P(\text{Jar}=\text{red}, \text{Fruit}=\text{apple})$$

$$= P(\text{Jar}=\text{red}) P(\text{Fruit}=\text{apple}|\text{Jar}=\text{red})$$

$$= 0.6 (0.167) = 0.1$$

*conditional
probability*

Doing all multiplications results in:



	apples	oranges	strawberries	
red jar (P=0.6)	0.6(0.167)	0.6(0.5)	0.6(0.33)	
	= 0.1	= 0.3	= 0.2	sum = 0.6

	apples	oranges	strawberries	
blue jar (P=0.4)	0.4(0.5)	0.4(0.167)	0.4(0.33)	
	= 0.2	= 0.067	= 0.133	sum = 0.4 ³

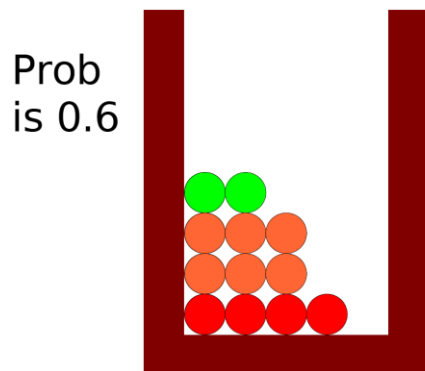
Joint Probability Table

- apple
- orange
- strawberry

Combine in a two-dimensional table to show joint probabilities of two events.

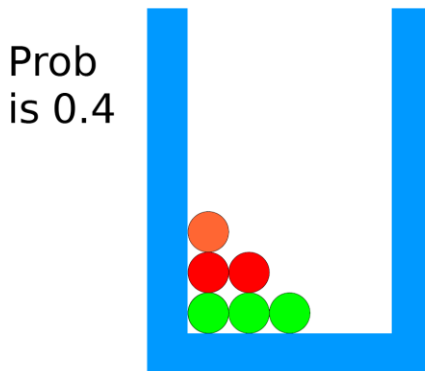
r = red, b = blue

a = apple, o = orange, s = strawberry



		Fruit			
		a	o	s	
Jar	r	0.1	0.3	0.2	$\Sigma = 0.6$
	b	0.2	0.067	0.133	$\Sigma = 0.4$
		$\Sigma = 0.3$	$\Sigma = 0.367$	$\Sigma = 0.333$	$\Sigma = 1.0$

Let J be random variable for Jar, and F be random variable for fruit.



		Fruit			
		a	o	s	
Jar	r	$P(J = r, F = a)$	$P(J = r, F = o)$	$P(J = r, F = s)$	$P(J = r)$
	b	$P(J = b, F = a)$	$P(J = b, F = o)$	$P(J = b, F = s)$	$P(J = b)$
		$P(F = a)$	$P(F = o)$	$P(F = s)$	1.0

Joint Probabilities and Bayes Rule

Just saw example of the *product rule*:

$$\begin{aligned} P(\text{Fruit}=\text{orange}, \text{Jar} = \text{blue}) \\ = P(\text{Fruit}=\text{orange} \mid \text{Jar} = \text{blue}) P(\text{Jar} = \text{blue}) \end{aligned}$$

Since $P(\text{Fruit}=\text{orange}, \text{Jar} = \text{blue}) = P(\text{Jar} = \text{blue}, \text{Fruit} = \text{orange})$,

$$\begin{aligned} P(\text{Jar} = \text{blue}, \text{Fruit}=\text{orange}) \\ = P(\text{Jar} = \text{blue} \mid \text{Fruit} = \text{orange}) P(\text{Fruit} = \text{orange}). \end{aligned}$$

Setting these equal leads to *Bayes Rule*:

$$\begin{aligned} P(\text{Jar} = \text{blue} \mid \text{Fruit} = \text{orange}) P(\text{Fruit} = \text{orange}) \\ = P(\text{Fruit}=\text{orange} \mid \text{Jar} = \text{blue}) P(\text{Jar} = \text{blue}) \end{aligned}$$

so

$$\begin{aligned} P(\text{Jar} = \text{blue} \mid \text{Fruit} = \text{orange}) \\ = P(\text{Fruit}=\text{orange} \mid \text{Jar} = \text{blue}) P(\text{Jar} = \text{blue}) / P(\text{Fruit} = \text{orange}) \end{aligned}$$

Joint Probabilities and Bayes Rule

On the right hand side of Bayes Rule, all terms are given to us except $P(\text{Fruit} = \text{orange})$:

$$\begin{aligned} P(\text{Jar} = \text{blue} \mid \text{Fruit} = \text{orange}) \\ = \frac{P(\text{Fruit}=\text{orange} \mid \text{Jar} = \text{blue}) P(\text{Jar} = \text{blue})}{P(\text{Fruit} = \text{orange})} \end{aligned}$$

We can use the *sum rule* to get this.

$$P(\text{Fruit} = \text{orange}) = \sum_j P(\text{Fruit}=\text{orange}, \text{Jar}=j) = 0.367$$

So, Bayes Rule can be rewritten as

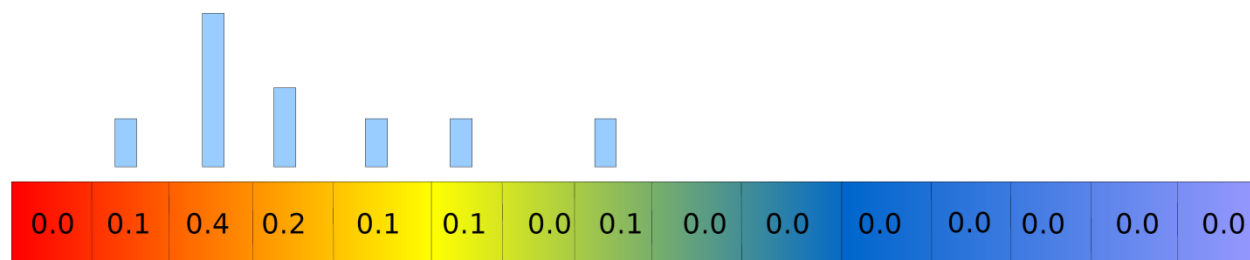
$$P(\text{Jar} = \text{blue} \mid \text{Fruit} = \text{orange})$$

$$\begin{aligned} &= \frac{P(\text{Fruit}=\text{orange} \mid \text{Jar} = \text{blue}) P(\text{Jar} = \text{blue})}{\sum_j P(\text{Fruit}=\text{orange}, \text{Jar}=j)} \\ &= \frac{P(\text{Fruit}=\text{orange} \mid \text{Jar} = \text{blue}) P(\text{Jar} = \text{blue})}{\sum_j P(\text{Fruit}=\text{orange} \mid \text{Jar}=j) P(\text{Jar} = j)} \end{aligned}$$

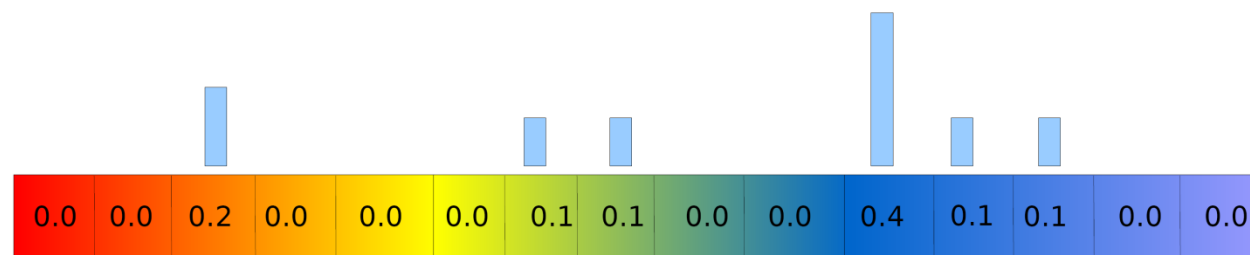
Probability Distributions

- Rather than three colors of fruit, imagine objects of 15 possible colors

Jar 1 contains objects with colors in these proportions



Jar 2 contains objects with colors in these proportions

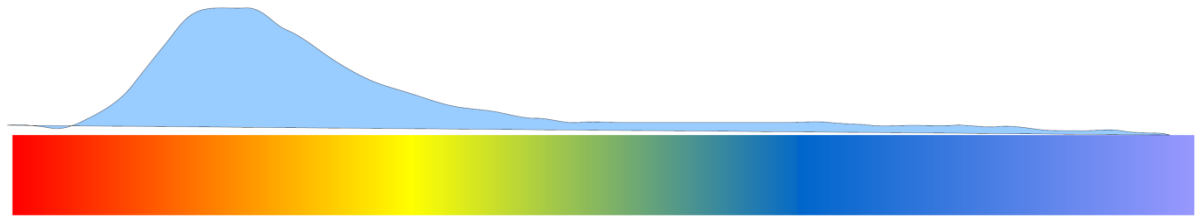


- Can calculate joint probability table as before.
- But what if we have 100 colors or 1000 colors?
- What if we have an infinite number of colors?

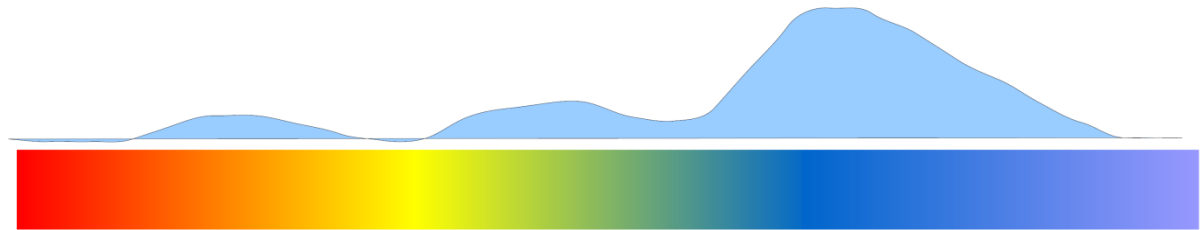
Distributions Over Continuous Values

- Probability of a color is a function over the continuous spectrum.

Jar 1 contains
objects with
colors in these
proportions



Jar 2 contains
objects with
colors in these
proportions

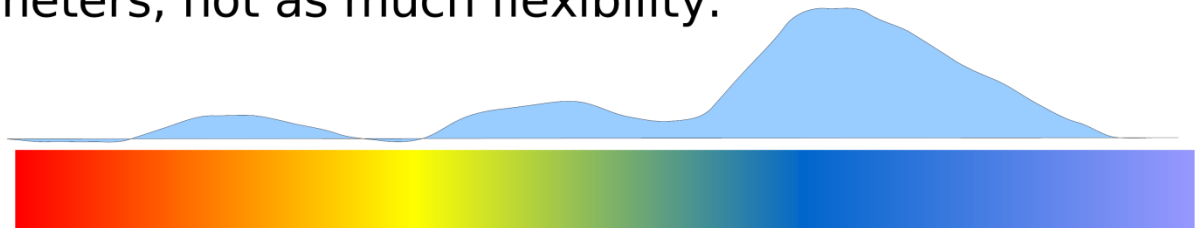


- But what function is this? Would require 1,000s of parameters to specify general function.
- Instead, let's use rather simple functions controlled by a few parameters.
- Common example: Gaussian (Normal) distribution

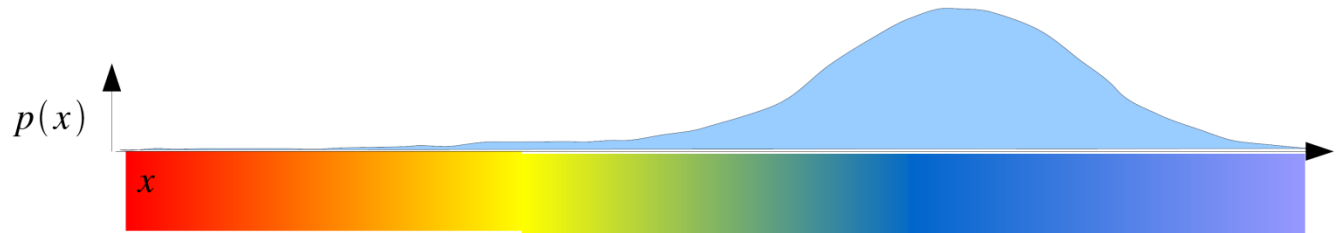
Gaussian Distribution

- With few parameters, not as much flexibility.

Jar 2 distribution



becomes



because

argument parameters

$$p(x; \mu, \sigma) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left(-\frac{1}{2\sigma^2}(x-\mu)^2\right)$$

- Easy to estimate parameters.

$$\mu = \frac{1}{N} \sum_{i=1}^N x_i \quad \sigma = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2$$

Gaussian Distribution

- Where do these expressions come from?

$$\mu = \frac{1}{N} \sum_{i=1}^N x_i \quad \sigma = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2$$

- Maximize the likelihood of the data.
 - Likelihood of data is product of probabilities of each sample x_i

$$p(X|\mu, \sigma) = \prod_{i=1}^N \left[\frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left(-\frac{1}{2\sigma^2}(x_i - \mu)^2\right) \right]$$

- Maximize this by maximizing its logarithm.

$$\ln p(X|\mu, \sigma) = -\frac{1}{2\sigma^2} \sum_{i=1}^N (x_i - \mu)^2 - \frac{N}{2} \ln \sigma^2 - \frac{N}{2} \ln(2\pi)$$

- Set its derivative with respect to μ to zero and solve for μ .
 - Set its derivative with respect to σ to zero and solve for σ .