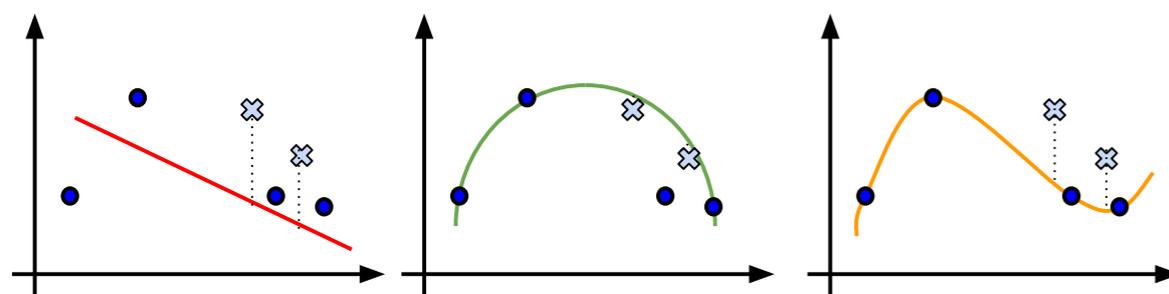
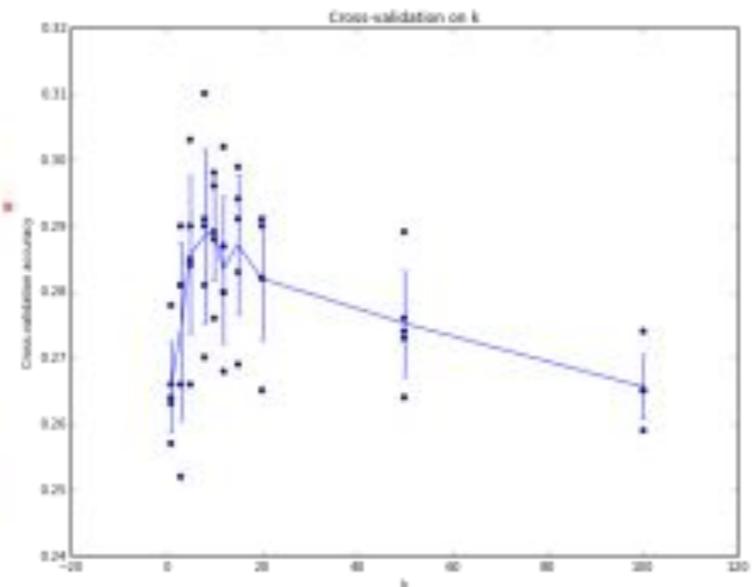
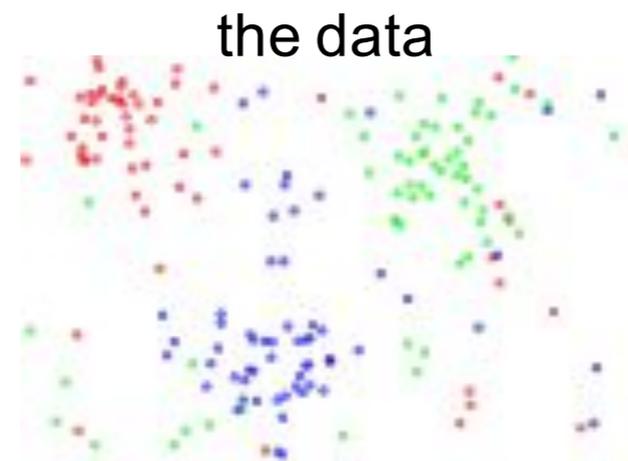
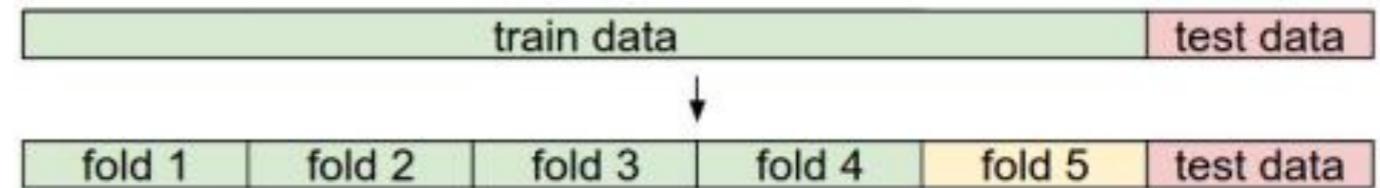
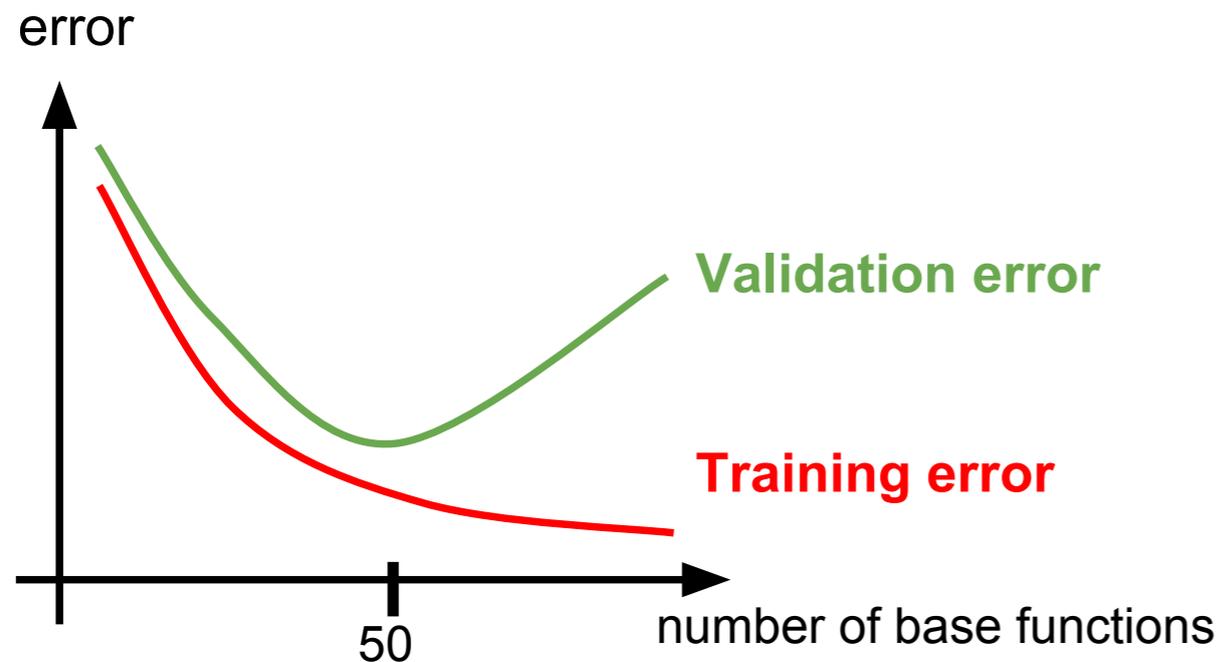


# BBM406

## Fundamentals of Machine Learning

Lecture 6:  
Learning theory  
Probability Review

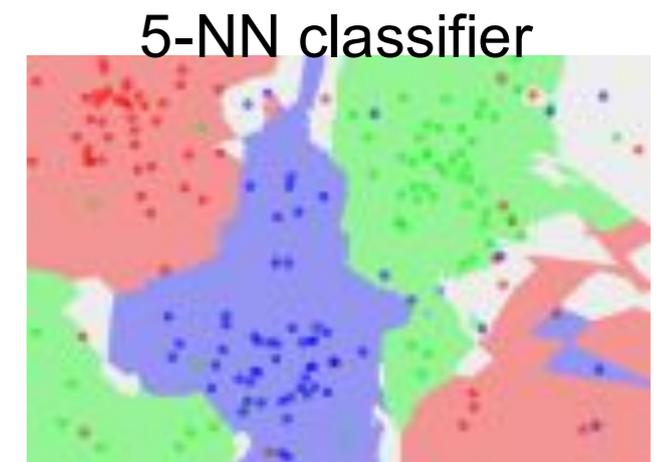
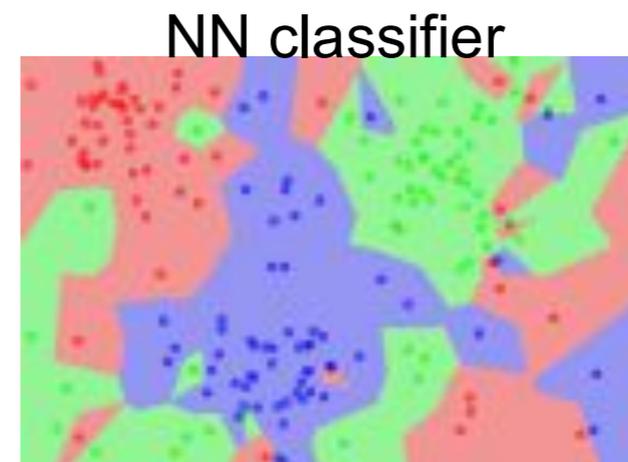
# Last time... Regularization, Cross-Validation



- Underfitting**
- large training error
  - large validation error

- Just Right**
- small training error
  - small validation error

- Overfitting**
- small training error
  - large validation error

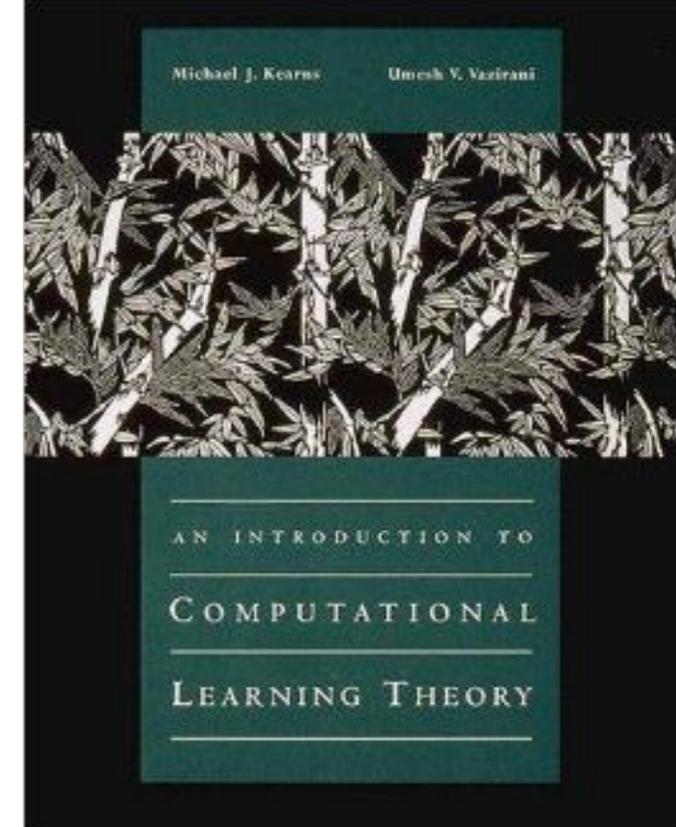


# Today

- Learning Theory
- Probability Review

# Learning Theory: Why ML Works

# Computational Learning Theory



- Entire subfield devoted to the mathematical analysis of machine learning algorithms
- Has led to several practical methods:
  - PAC (probably approximately correct) learning
    - boosting
  - VC (Vapnik–Chervonenkis) theory
    - support vector machines

Annual conference: [Conference on Learning Theory \(COLT\)](#)

# The Role of Theory

- Theory can serve two roles:
  - It can justify and help understand why common practice works.  
*theory after*
  - It can also serve to suggest new algorithms and approaches that turn out to work well in practice.  
*theory before*

**Often, it turns out to be a mix!**

# The Role of Theory

- Practitioners discover something that works surprisingly well.
- Theorists figure out why it works and prove something about it.
  - In the process, they make it better or find new algorithms.
- Theory can also help you understand **what's possible and what's not possible.**

# Learning and Inference

The inductive inference process:

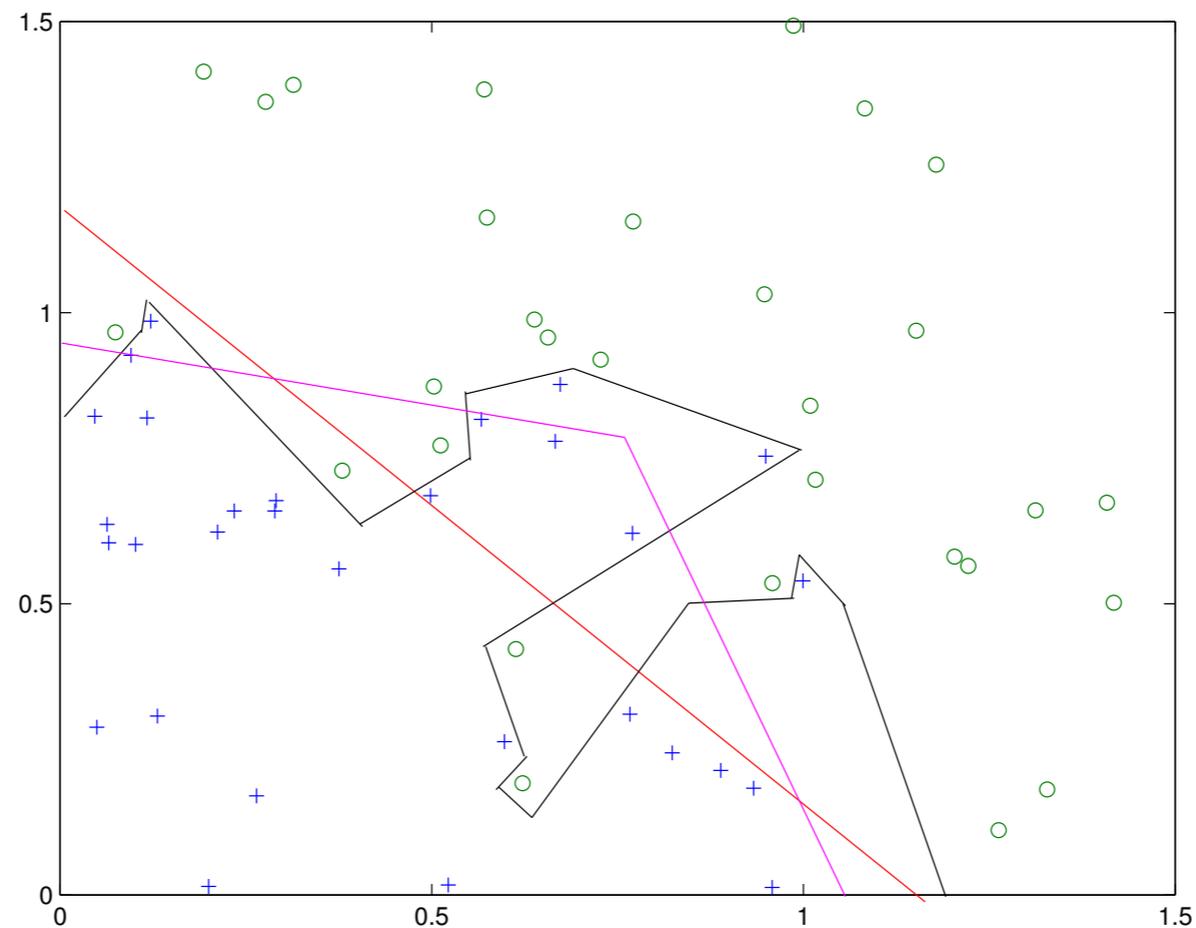
1. Observe a phenomenon
  2. Construct a model of the phenomenon
  3. Make predictions
- This is more or less the definition of natural sciences !
  - The goal of Machine Learning is to **automate** this process
  - The goal of Learning Theory is to **formalize** it.

# Pattern recognition

- We consider here the **supervised learning** framework for pattern recognition:
  - Data consists of pairs (instance, label)
  - Label is +1 or -1
  - Algorithm constructs a function (instance  $\rightarrow$  label)
  - Goal: make few mistakes on future unseen instances

# Approximation/Interpolation

- It is always possible to build a function that fits exactly the data.



- But is it reasonable?

# Occam's Razor

- Idea: look for **regularities** in the observed phenomenon

These can be **generalized** from the observed past to the future

⇒ choose the simplest consistent model

- How to measure simplicity ?
  - Physics: number of constants
  - Description length
  - Number of parameters
  - ...



William of Occam  
(c. 1288 – c. 1348)

# No Free Lunch

- **No Free Lunch**
  - if there is no assumption on how the **past** is related to the future, prediction is impossible
  - if there is no **restriction** on the possible phenomena, generalization is impossible
- We need to make assumptions
- Simplicity is not absolute
- Data will never replace knowledge
- Generalization = data + knowledge

# Probably Approximately Correct (PAC) Learning

- A formalism based on the realization that the best we can hope of an algorithm is that
  - It does a good job most of the time (**probably approximately correct**)

# Probably Approximately Correct (PAC) Learning

- Consider a hypothetical learning algorithm
  - We have 10 different binary classification data sets.
  - For each one, it comes back with functions  $f_1, f_2, \dots, f_{10}$ .
  - ♦ For some reason, whenever you run  $f_4$  on a test point, it crashes your computer. For the other learned functions, their performance on test data is always at most 5% error.
  - ♦ If this situation is guaranteed to happen, then this hypothetical learning algorithm is a PAC learning algorithm.
  - ❖ It satisfies **probably** because it only failed in one out of ten cases, and it's **approximate** because it achieved low, but non-zero, error on the remainder of the cases.

# PAC Learning

**Definitions 1.** *An algorithm  $A$  is an  $(\epsilon, \delta)$ -PAC learning algorithm if, for all distributions  $\mathcal{D}$ : given samples from  $\mathcal{D}$ , the probability that it returns a “bad function” is at most  $\delta$ ; where a “bad” function is one with test error rate more than  $\epsilon$  on  $\mathcal{D}$ .*

# PAC Learning

- Two notions of efficiency
  - **Computational complexity:** Prefer an algorithm that runs quickly to one that takes forever
  - **Sample complexity:** The number of examples required for your algorithm to achieve its goals

**Definition:** An algorithm  $\mathcal{A}$  is an **efficient  $(\epsilon, \delta)$ -PAC learning algorithm** if it is an  $(\epsilon, \delta)$ -PAC learning algorithm whose runtime is polynomial in  $\frac{1}{\epsilon}$  and  $\frac{1}{\delta}$ .

*In other words, to let your algorithm to achieve 4% error rather than 5%, the runtime required to do so should not go up by an exponential factor!*

# Example: PAC Learning of Conjunctions

- Data points are binary vectors, for instance  $\mathbf{x} = \langle 0, 1, 1, 0, 1 \rangle$
- Some Boolean conjunction defines the true labeling of this data (e.g.  $x_1 \wedge x_2 \wedge x_5$ )
- There is some distribution  $\mathcal{D}_X$  over binary data points (vectors)  
 $\mathbf{x} = \langle x_1, x_2, \dots, x_D \rangle$ .
- There is a fixed concept conjunction  $c$  that we are trying to learn.
- There is no noise, so for any example  $x$ , its true label is simply  $y = c(\mathbf{x})$

- **Example:**

- Clearly, the true formula cannot include the terms  $x_1, x_2, \neg x_3, \neg x_4$

$y$	$x_1$	$x_2$	$x_3$	$x_4$
+1	0	0	1	1
+1	0	1	1	1
-1	1	1	0	1

# Example: PAC Learning of Conjunctions

$y$	$x_1$	$x_2$	$x_3$	$x_4$
+1	0	0	1	1
+1	0	1	1	1
-1	1	1	0	1

$$f^0(\mathbf{x}) = x_1 \wedge \neg x_1 \wedge x_2 \wedge \neg x_2 \wedge x_3 \wedge \neg x_3 \wedge x_4 \wedge \neg x_4$$

$$f^1(\mathbf{x}) = \neg x_1 \wedge \neg x_2 \wedge x_3 \wedge x_4$$

$$f^2(\mathbf{x}) = \neg x_1 \wedge x_3 \wedge x_4$$

$$f^3(\mathbf{x}) = \neg x_1 \wedge x_3 \wedge x_4$$

- After processing an example, it is **guaranteed to classify that example correctly** (provided that there is no noise)
- **Computationally very efficient**
  - Given a data set of  $N$  examples in  $D$  dimensions, it takes  $O(ND)$  time to process the data. This is linear in the size of the data set.

## Algorithm 30 BINARYCONJUNCTIONTRAIN(D)

```

1:  $f \leftarrow x_1 \wedge \neg x_1 \wedge x_2 \wedge \neg x_2 \wedge \dots \wedge x_D \wedge \neg x_D$  // initialize function
2: for all positive examples  $(x_{,+1})$  in  $D$  do
3:   for  $d = 1 \dots D$  do
4:     if  $x_d = 0$  then
5:        $f \leftarrow f$  without term " $x_d$ "
6:     else
7:        $f \leftarrow f$  without term " $\neg x_d$ "
8:     end if
9:   end for
10: end for
11: return  $f$ 

```

“Throw Out Bad Terms”

# Example: PAC Learning of Conjunctions

$y$	$x_1$	$x_2$	$x_3$	$x_4$
+1	0	0	1	1
+1	0	1	1	1
-1	1	1	0	1

## Algorithm 30 BINARYCONJUNCTIONTRAIN( $\mathbf{D}$ )

```

1:  $f \leftarrow x_1 \wedge \neg x_1 \wedge x_2 \wedge \neg x_2 \wedge \dots \wedge x_D \wedge \neg x_D$  // initialize function
2: for all positive examples  $(x, +1)$  in  $\mathbf{D}$  do
3:   for  $d = 1 \dots D$  do
4:     if  $x_d = 0$  then
5:        $f \leftarrow f$  without term " $x_d$ "
6:     else
7:        $f \leftarrow f$  without term " $\neg x_d$ "
8:     end if
9:   end for
10: end for
11: return  $f$ 

```

“Throw Out Bad Terms”

- Is this an efficient  $(\epsilon, \delta)$ -PAC learning algorithm?
- What about **sample complexity**?
  - How many examples  $N$  do you need to see in order to guarantee that it achieves an error rate of at most  $\epsilon$  (in all but  $\delta$ - many cases)?
  - Perhaps  $N$  has to be gigantic (like  $2^{2^{D/\epsilon}}$ ) to (probably) guarantee a small error.

# Vapnik-Chervonenkis (VC) Dimension

- A classic measure of complexity of infinite hypothesis classes based on this intuition.
- The VC dimension is a very classification-oriented notion of complexity
  - The idea is to look at a finite set of unlabeled examples
  - no matter how these points were labeled, would we be able to find a hypothesis that correctly classifies them
- The idea is that as you add more points, being able to represent an arbitrary labeling becomes harder and harder.

**Definitions 2.** For data drawn from some space  $\mathcal{X}$ , the *VC dimension* of a hypothesis space  $\mathcal{H}$  over  $\mathcal{X}$  is the maximal  $K$  such that: *there exists a set  $X \subseteq \mathcal{X}$  of size  $|X| = K$ , such that **for all** binary labelings of  $X$ , there exists a function  $f \in \mathcal{H}$  that matches this labeling.*

# How many points can a linear boundary classify exactly? (1-D)

- 2 points:

Yes!



- 3 points:

No!



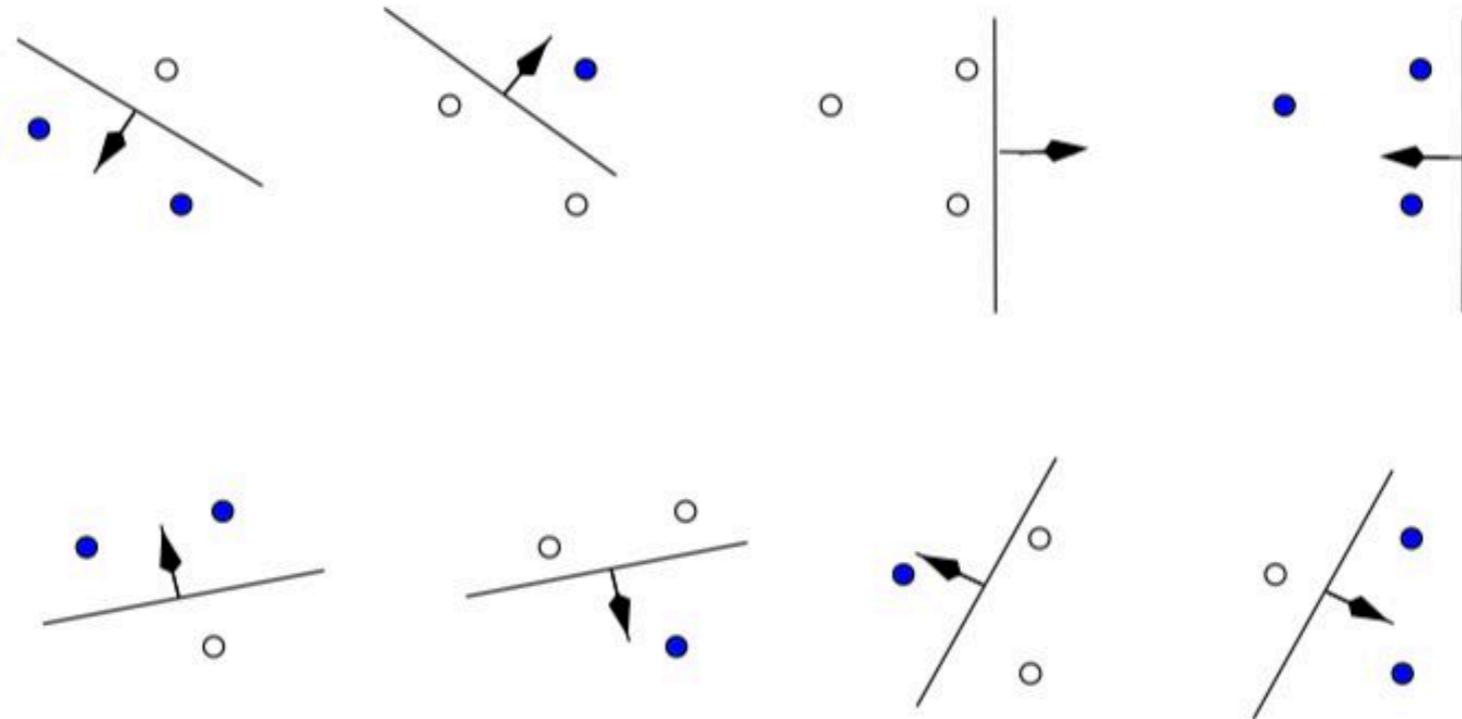
etc (8 total)

**VC-dimension = 2**

# How many points can a linear boundary classify exactly? (2-D)

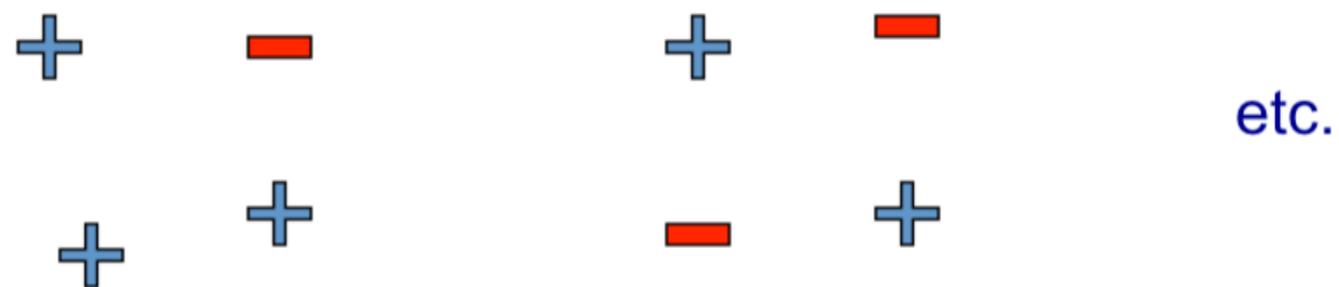
- 3 points:

Yes!



- 4 points:

No!



**VC-dimension = 3**

# Basic Probability Review

# Probability

- A is non-deterministic event
  - Can think of A as a boolean-valued variable
- Examples
  - A = your next patient has cancer
  - A = Rafael Nadal wins French Open 2019

# Interpreting Probabilities



If I flip this coin, the probability that it will come up heads is 0.5

- **Frequentist Interpretation:** If we flip this coin many times, it will come up heads about half the time. *Probabilities are the expected frequencies of events over repeated trials.*
- **Bayesian Interpretation:** I believe that my next toss of this coin is equally likely to come up heads or tails. *Probabilities quantify subjective beliefs about single events.*
- Viewpoints play complementary roles in **machine learning:**
  - Bayesian view used to build models based on domain knowledge, and automatically derive learning algorithms
  - Frequentist view used to analyze worst case behavior of learning algorithms, in limit of large datasets
- From either view, basic mathematics is the same!



# Axioms of Probability

- $0 \leq P(A) \leq 1$
- $P(\text{empty-set}) = 0$
- $P(\text{everything}) = 1$
- $P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$

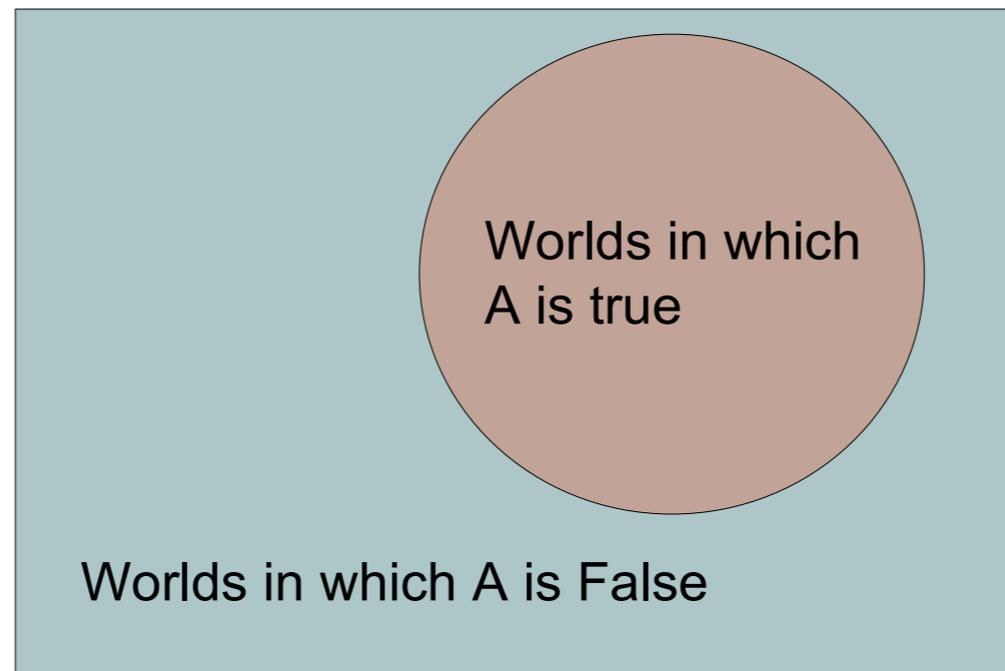
# Interpreting the Axioms

- $0 \leq P(A) \leq 1$
- $P(\text{empty-set}) = 0$
- $P(\text{everything}) = 1$
- $P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$

Event space of  
all possible  
worlds



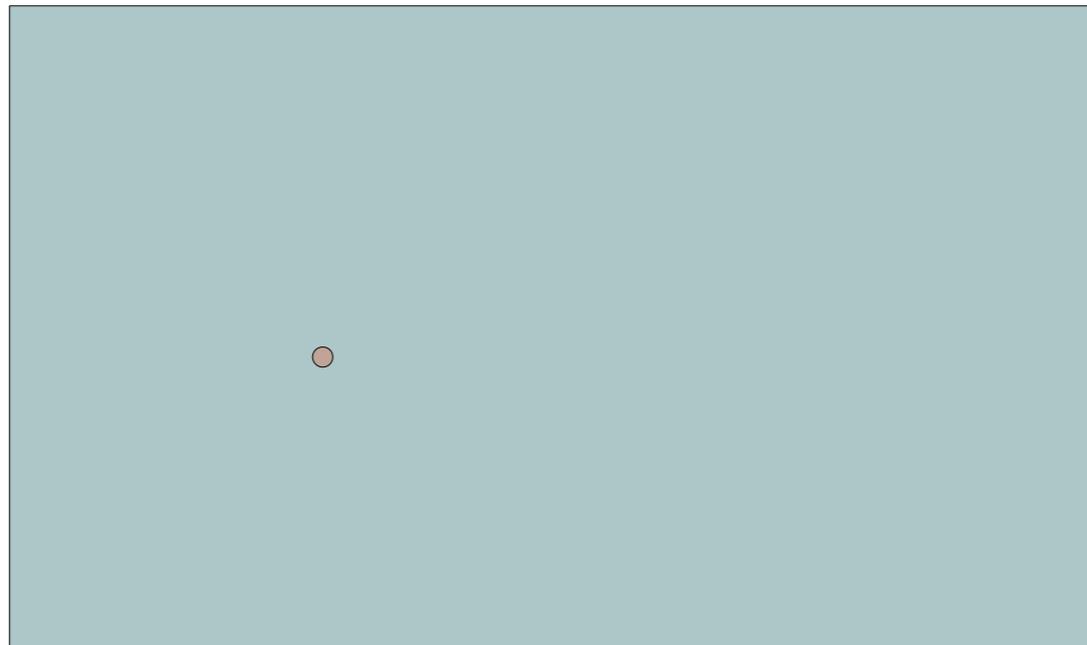
Its area is 1



$P(A) = \text{Area of reddish oval}$

# Interpreting the Axioms

- $0 \leq P(A) \leq 1$
- $P(\text{empty-set}) = 0$
- $P(\text{everything}) = 1$
- $P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$

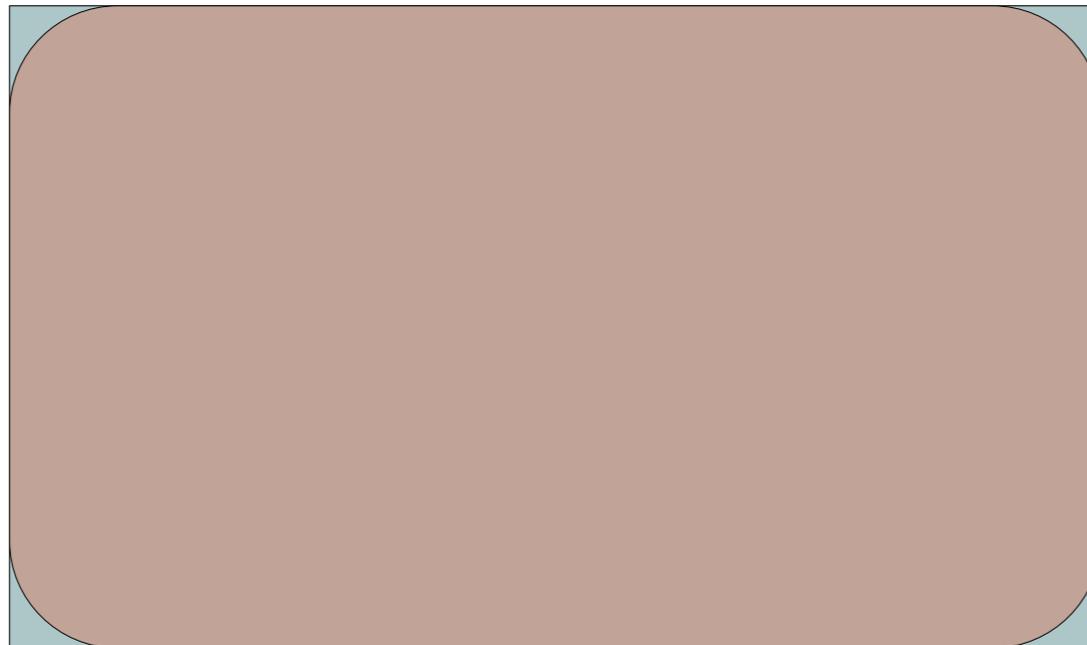


The area of A can't get any smaller than 0

And a zero area would mean no world could ever have A true

# Interpreting the Axioms

- $0 \leq P(A) \leq 1$
- $P(\text{empty-set}) = 0$
- $P(\text{everything}) = 1$
- $P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$

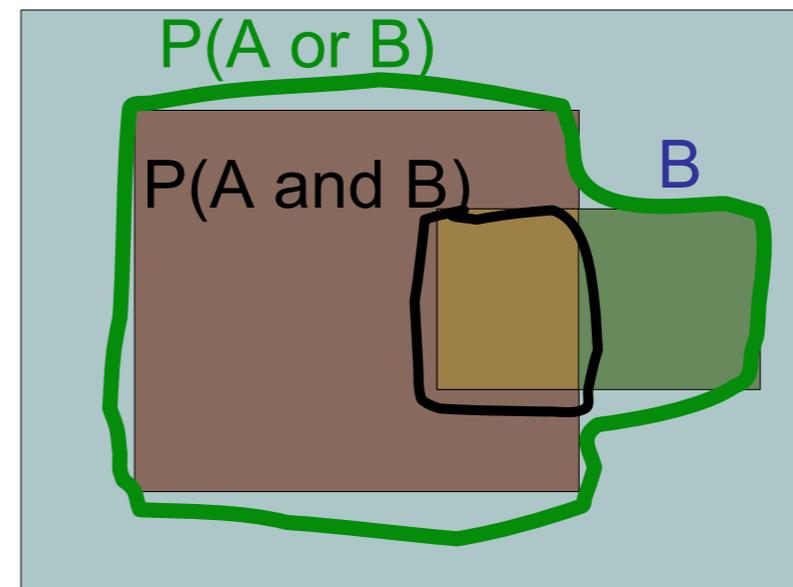
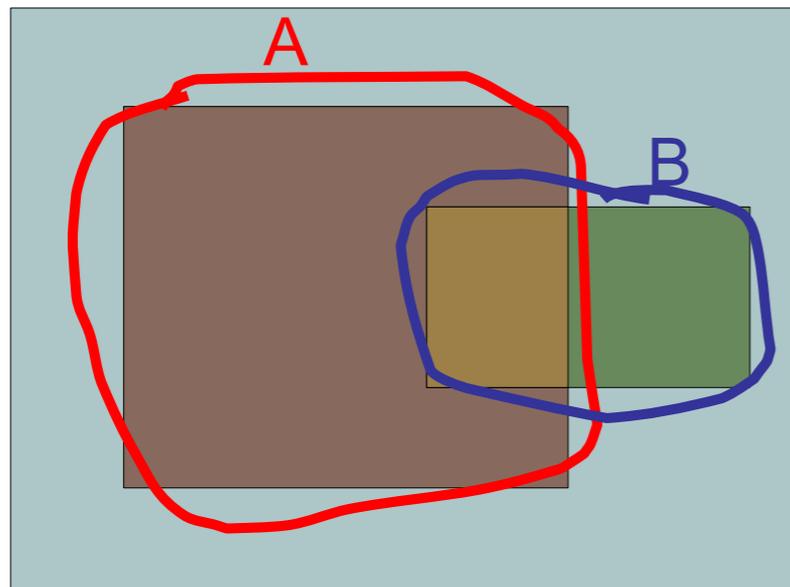


The area of A can't get any bigger than 1

And an area of 1 would mean all worlds will have A true

# Interpreting the Axioms

- $0 \leq P(A) \leq 1$
- $P(\text{empty-set}) = 0$
- $P(\text{everything}) = 1$
- $P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$



Simple addition and subtraction

# Discrete Random Variables

$X$   $\longrightarrow$  discrete random variable

$\mathcal{X}$   $\longrightarrow$  sample space of possible outcomes,  
which may be finite or countably infinite

$x \in \mathcal{X}$   $\longrightarrow$  outcome of sample of discrete random variable

# Discrete Random Variables

$X$   $\longrightarrow$  discrete random variable

$\mathcal{X}$   $\longrightarrow$  sample space of possible outcomes,  
which may be finite or countably infinite

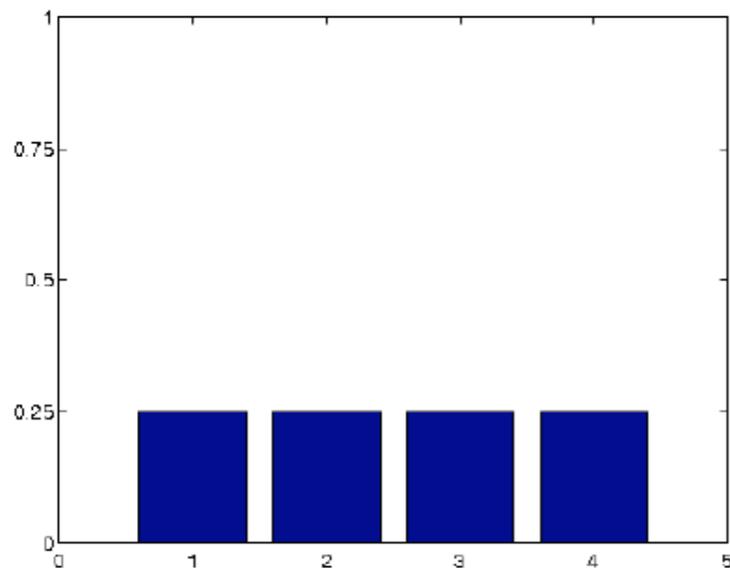
$x \in \mathcal{X}$   $\longrightarrow$  outcome of sample of discrete random variable

$p(X = x)$   $\longrightarrow$  probability distribution (probability mass function)

$p(x)$   $\longrightarrow$  shorthand used when no ambiguity

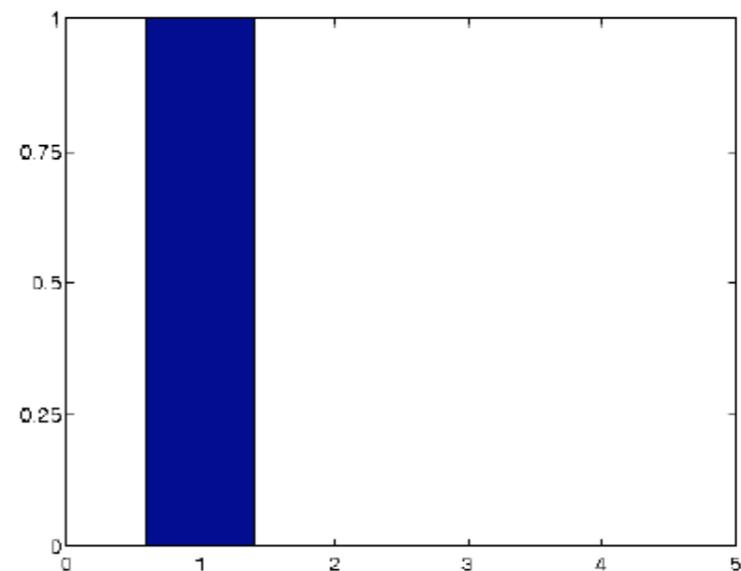
$0 \leq p(x) \leq 1$  for all  $x \in \mathcal{X}$

$$\sum_{x \in \mathcal{X}} p(x) = 1$$



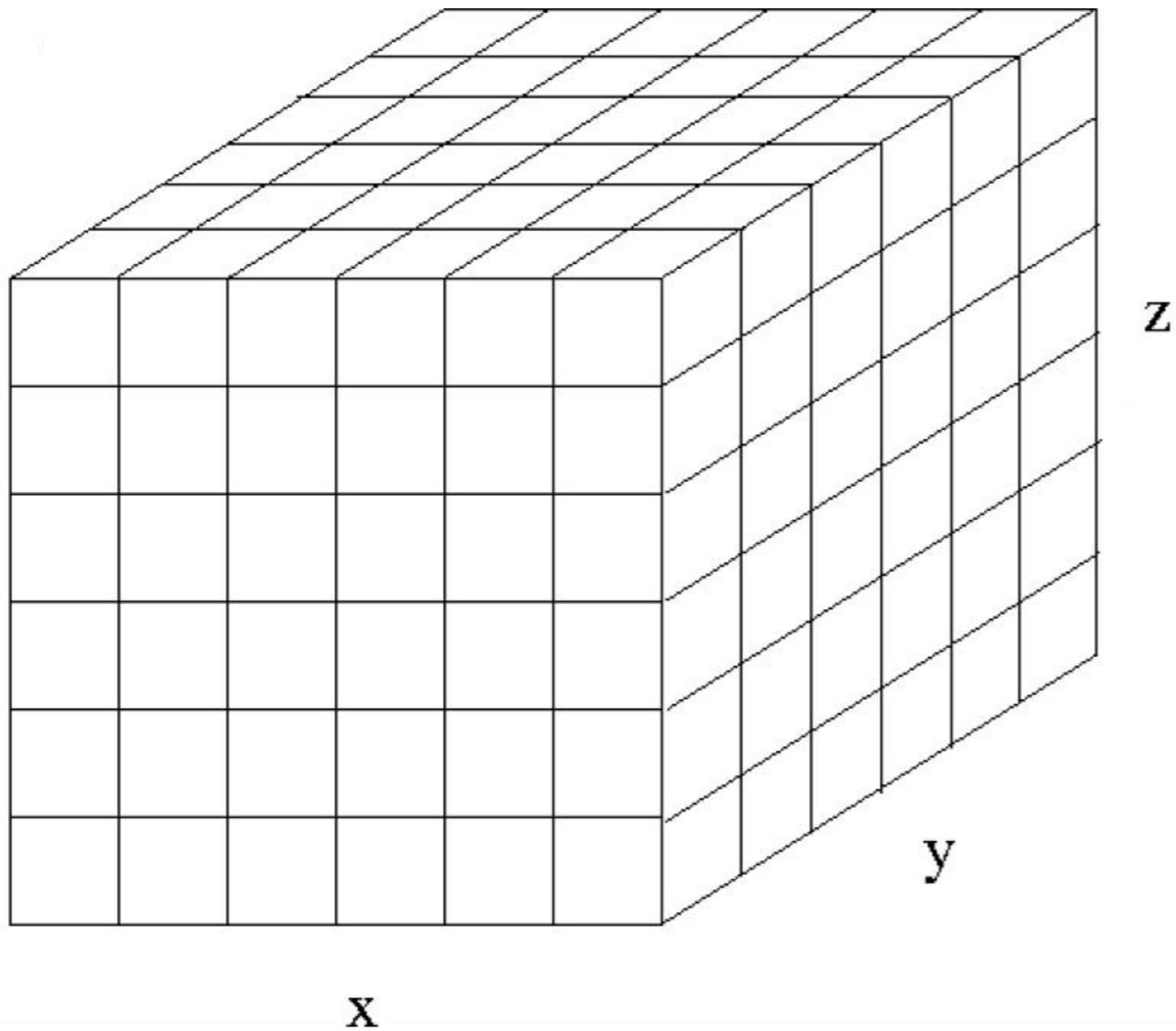
*uniform distribution*

$$\mathcal{X} = \{1, 2, 3, 4\}$$



*degenerate distribution*

# Joint Distribution



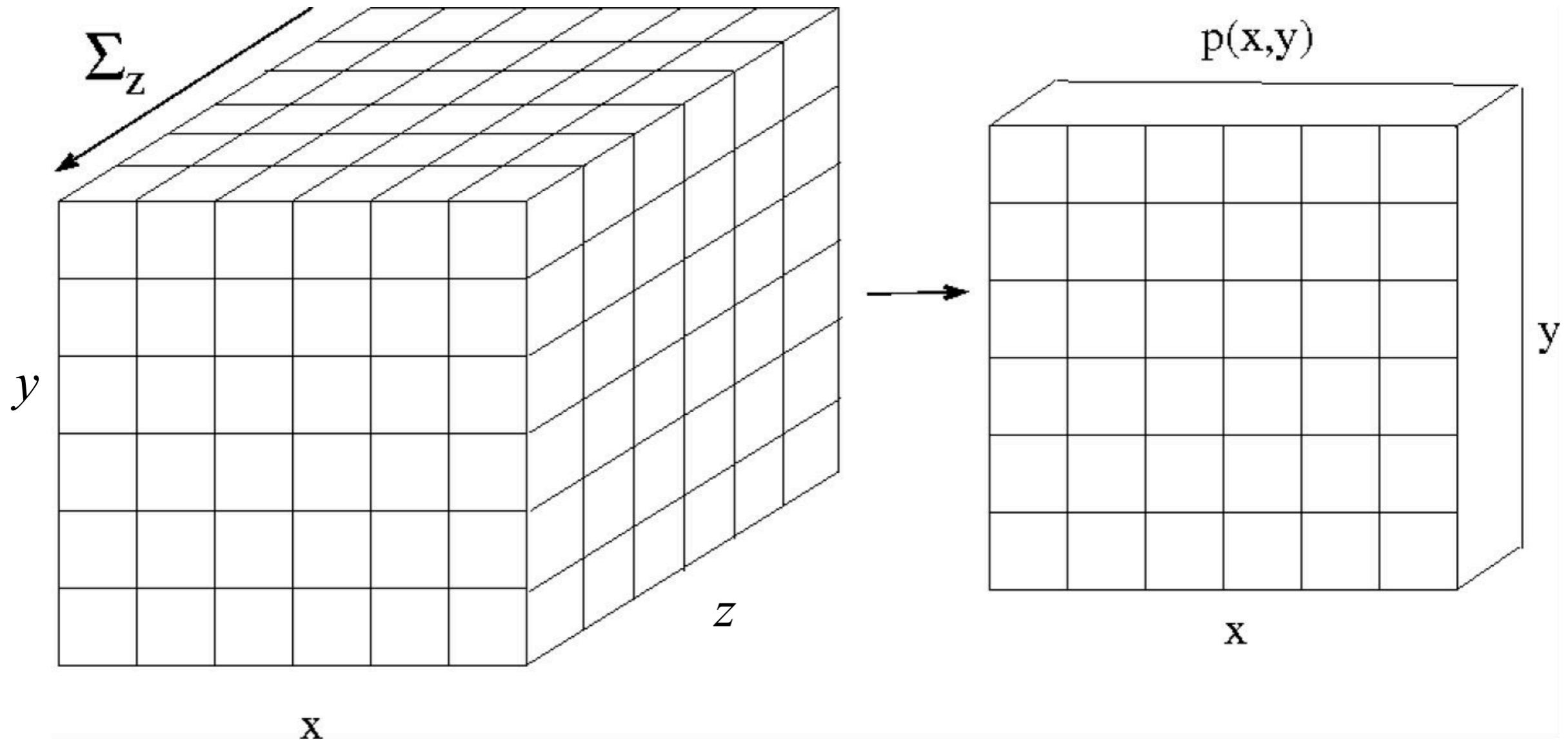
# Marginalization

- Marginalization

- Events:  $P(A) = P(A \text{ and } B) + P(A \text{ and not } B)$

- Random variables  $P(X = x) = \sum_y P(X = x, Y = y)$

# Marginal Distributions



$$p(x, y) = \sum_{z \in \mathcal{Z}} p(x, y, z)$$

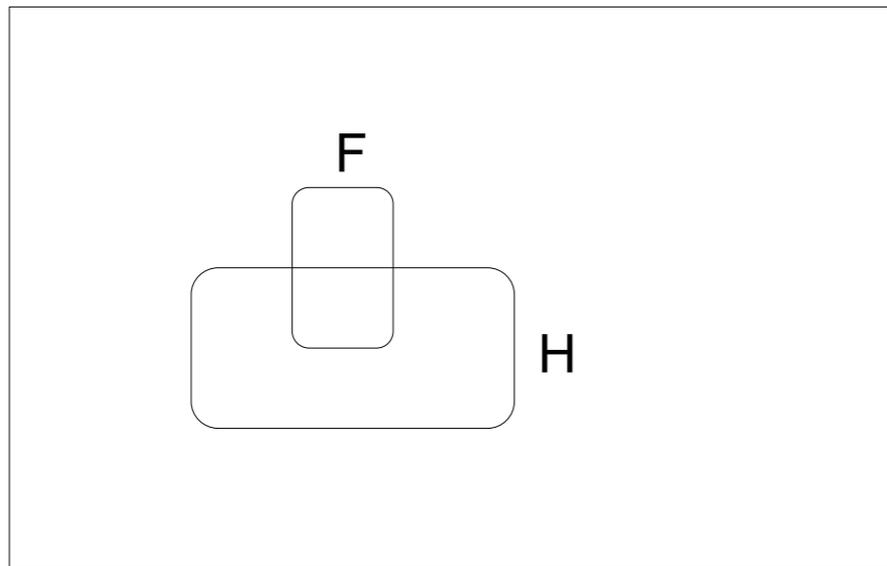
$$p(x) = \sum_{y \in \mathcal{Y}} p(x, y)$$

# Conditional Probabilities

- $P(Y=y \mid X=x)$
- What do you believe about  $Y=y$ , if I tell you  $X=x$ ?
- $P(\text{Rafael Nadal wins French Open 2019})?$
- What if I tell you:
  - He has won the French Open 11/13 he has played there
  - Rafael Nadal is ranked 1

# Conditional Probabilities

- $P(A | B)$  = In worlds that where B is true, fraction where A is true
- Example
  - H: “Have a headache”
  - F: “Coming down with Flu”



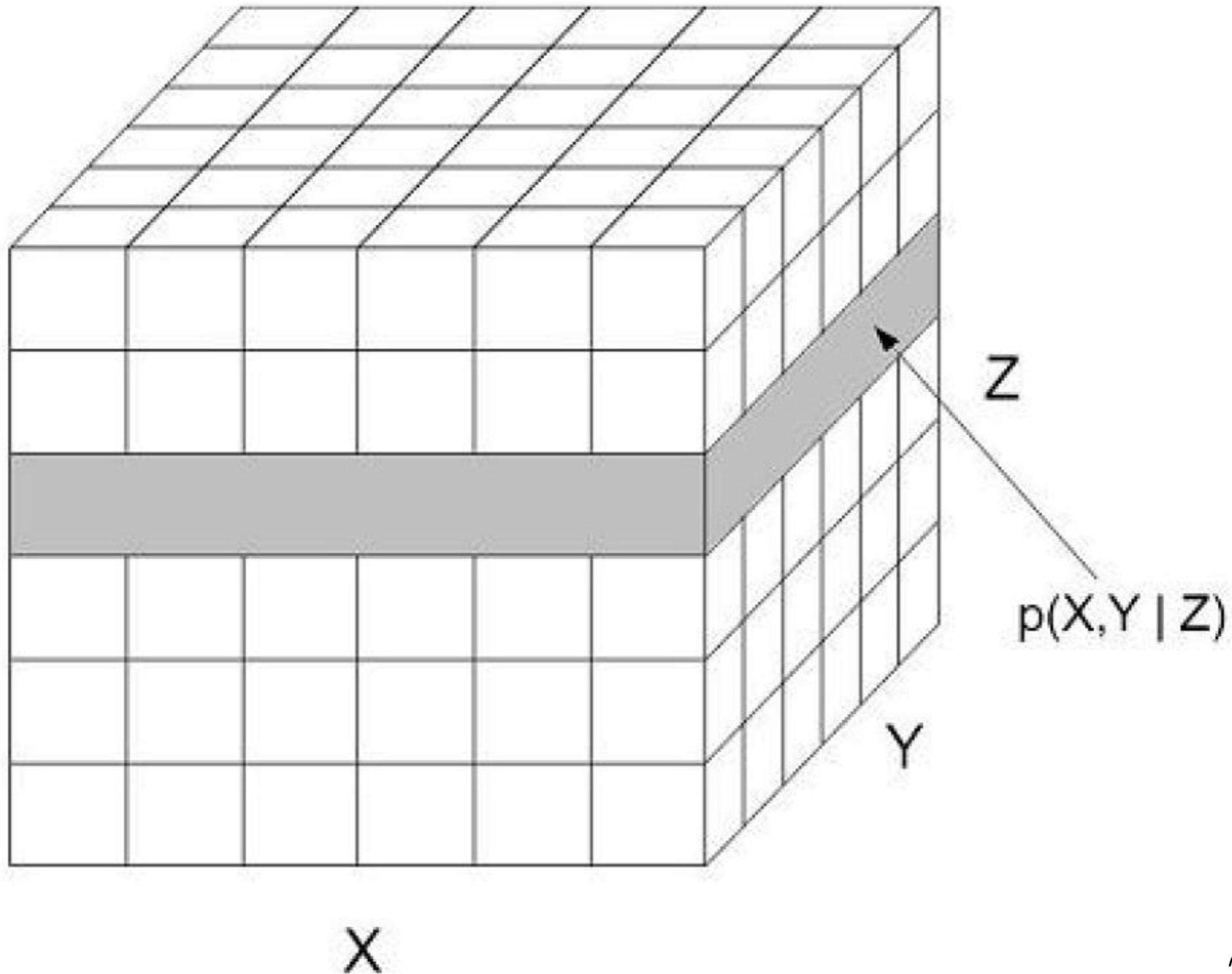
$$P(H) = 1/10$$

$$P(F) = 1/40$$

$$P(H|F) = 1/2$$

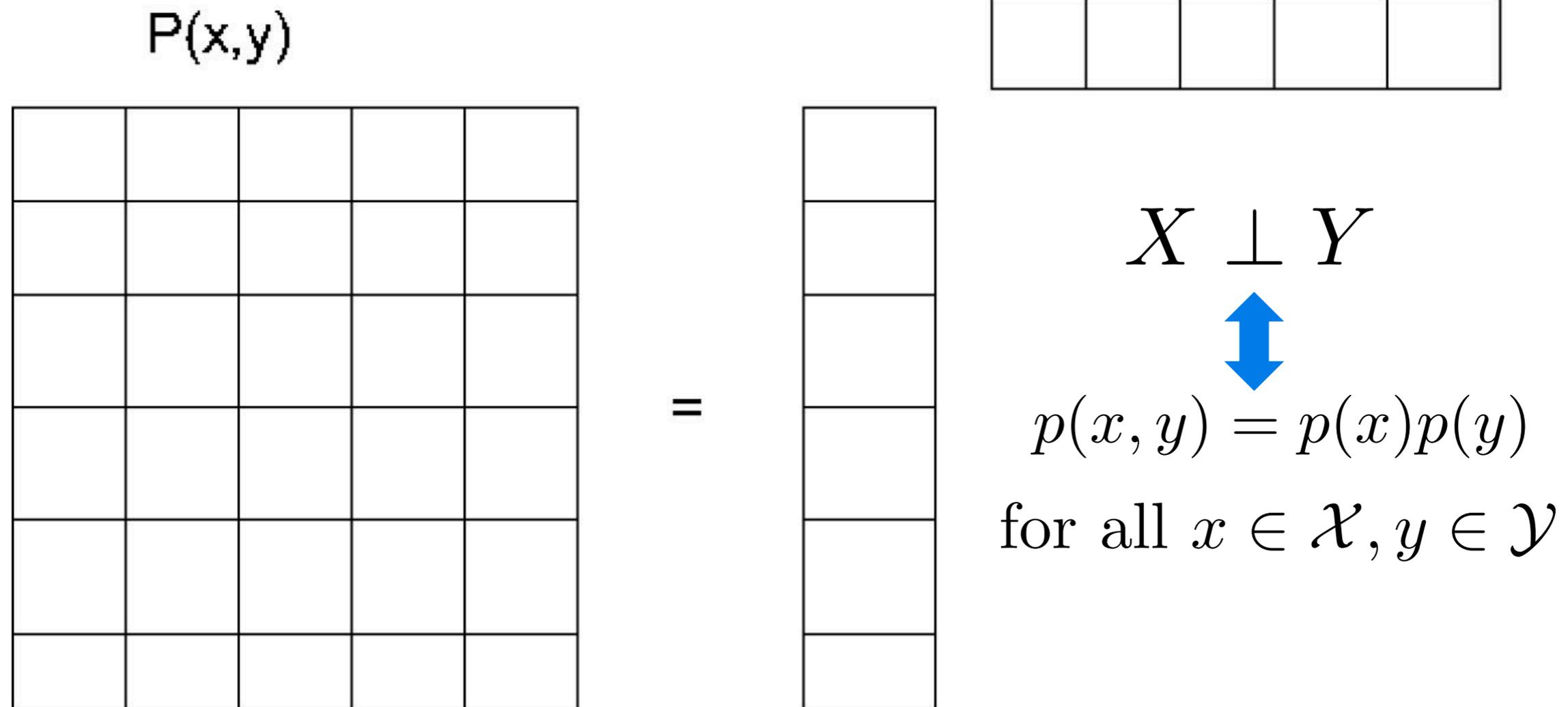
Headaches are rare and flu is rarer, but if you're coming down with flu there's a 50-50 chance you'll have a headache.

# Conditional Distributions



$$p(x, y | Z = z) = \frac{p(x, y, z)}{p(z)}$$

# Independent Random Variables



Equivalent conditions on conditional probabilities:

$$p(x \mid Y = y) = p(x) \text{ and } p(y) > 0 \text{ for all } y \in \mathcal{Y}$$

$$p(y \mid X = x) = p(y) \text{ and } p(x) > 0 \text{ for all } x \in \mathcal{X}$$

# Bayes Rule (Bayes Theorem)

$$p(x, y) = p(x)p(y | x) = p(y)p(x | y)$$

$$p(y | x) = \frac{p(x, y)}{p(x)} = \frac{p(x | y)p(y)}{\sum_{y' \in \mathcal{Y}} p(y')p(x | y')} \\ \propto p(x | y)p(y)$$



- A basic identity from the definition of conditional probability
- Used in ways that have no thing to do with Bayesian statistics!
- Typical application to learning and data analysis:

$Y$	$\longrightarrow$	unknown parameters we would like to infer
$X = x$	$\longrightarrow$	observed data available for learning
$p(y)$	$\longrightarrow$	prior distribution (domain knowledge)
$p(x   y)$	$\longrightarrow$	likelihood function (measurement model)
$p(y   x)$	$\longrightarrow$	posterior distribution (learned information)

# Binary Random Variables

- **Bernoulli Distribution:** Single toss of a (possibly biased) coin

$$\mathcal{X} = \{0, 1\}$$

$$0 \leq \theta \leq 1$$

$$\text{Ber}(x \mid \theta) = \theta^{\delta(x,1)} (1 - \theta)^{\delta(x,0)}$$



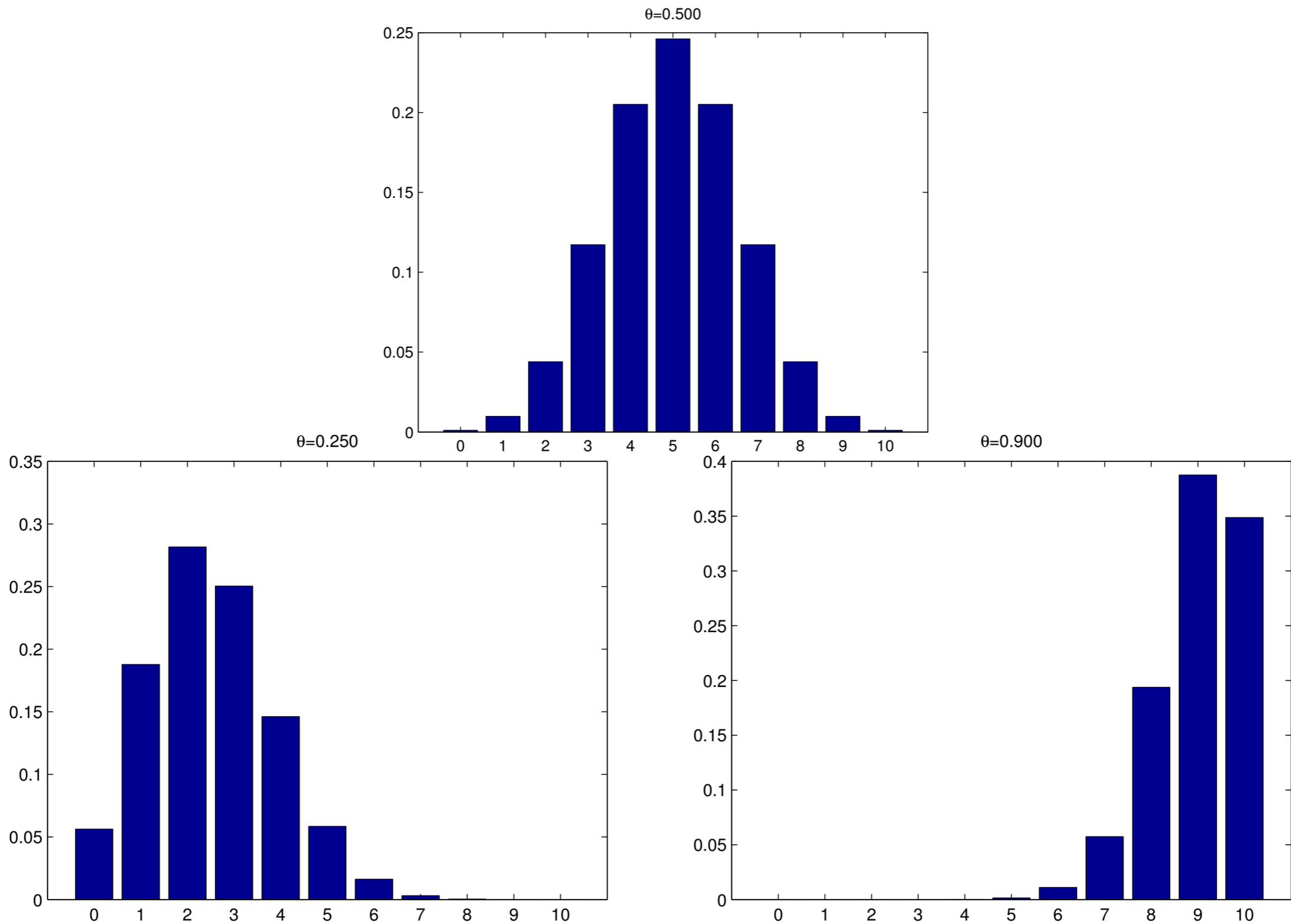
- **Binomial Distribution:** Toss a single (possibly biased) coin  $n$  times, and report the number  $k$  of times it comes up

$$\mathcal{K} = \{0, 1, 2, \dots, n\}$$

$$0 \leq \theta \leq 1$$

$$\text{Bin}(k \mid n, \theta) = \binom{n}{k} \theta^k (1 - \theta)^{n-k} \quad \binom{n}{k} = \frac{n!}{(n-k)!k!}$$

# Binomial Distributions



# Bean Machine (Sir Francis Galton)



[http://en.wikipedia.org/wiki/Bean\\_machine](http://en.wikipedia.org/wiki/Bean_machine)

# Categorical Random Variables

- **Multinoulli Distribution:** Single roll of a (possibly biased) die

$$\mathcal{X} = \{0, 1\}^K, \sum_{k=1}^K x_k = 1 \quad \text{binary vector encoding}$$

$$\theta = (\theta_1, \theta_2, \dots, \theta_K), \theta_k \geq 0, \sum_{k=1}^K \theta_k = 1$$

$$\text{Cat}(x | \theta) = \prod_{k=1}^K \theta_k^{x_k}$$

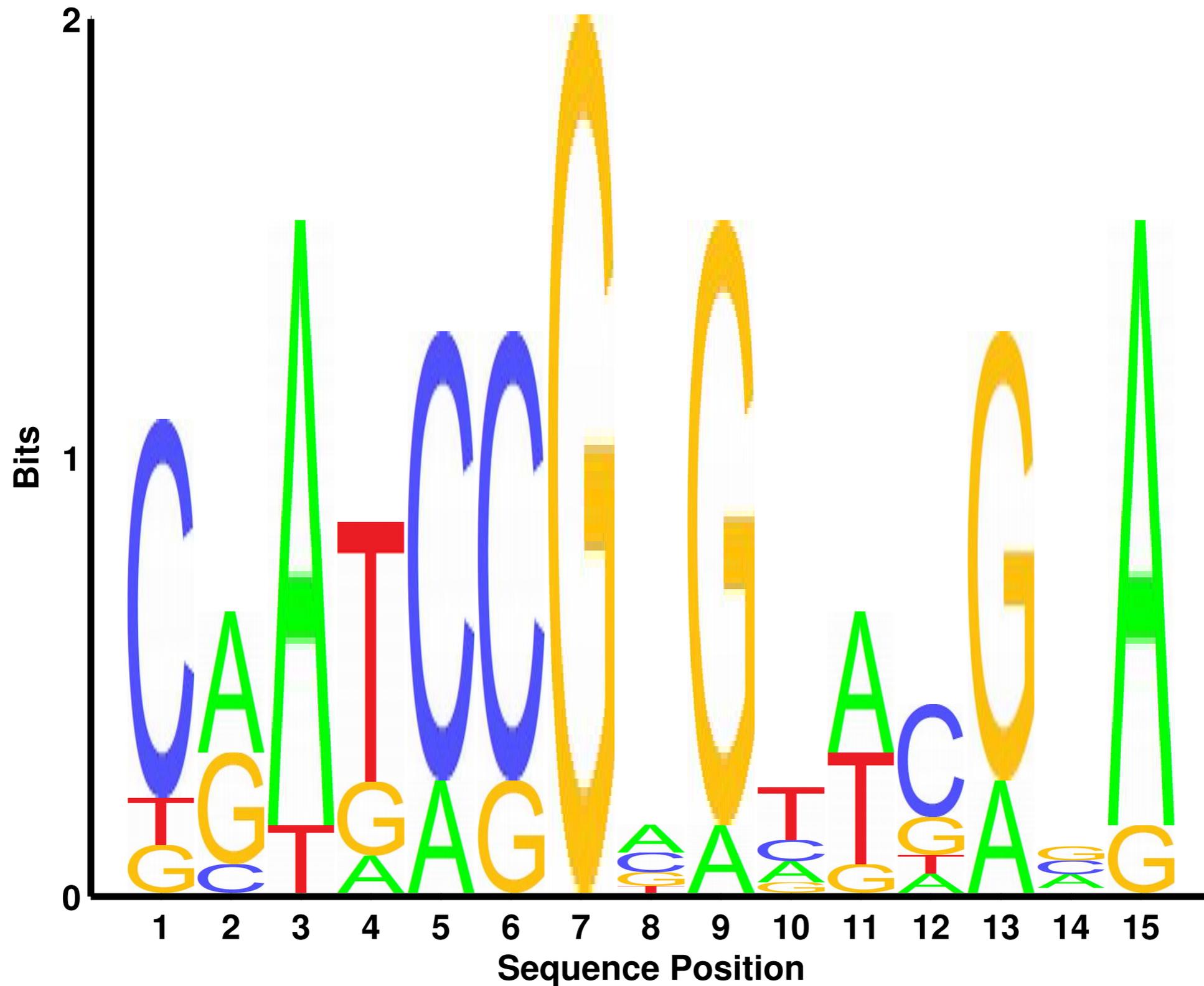
- **Multinomial Distribution:** Roll a single (possibly biased) die  $n$  times, and report the number  $n_k$  of each possible outcome

$$\text{Mu}(x | n, \theta) = \binom{n}{n_1 \dots n_K} \prod_{k=1}^K \theta_k^{n_k} \quad n_k = \sum_{i=1}^n x_{ik}$$

# Aligned DNA Sequences

```
c g a t a c g g g g t c g a a  
c a a t c c g a g a t c g c a  
c a a t c c g t g t t g g g a  
c a a t c g g c a t g c g g g  
c g a g c c g c g t a c g a a  
c a t a c g g a g c a c g a a  
t a a t c c g g g c a t g t a  
c g a g c c g a g t a c a g a  
c c a t c c g c g t a a g c a  
g g a t a c g a g a t g a c a
```

# Multinomial Model of DNA



**Next Lecture:**  
**Maximum Likelihood Estimation**  
**(MLE)**