

# Machine Learning 2

DS 4420 / ML 2

## Math review

Byron C Wallace



# Probability

# Examples: Independent Events

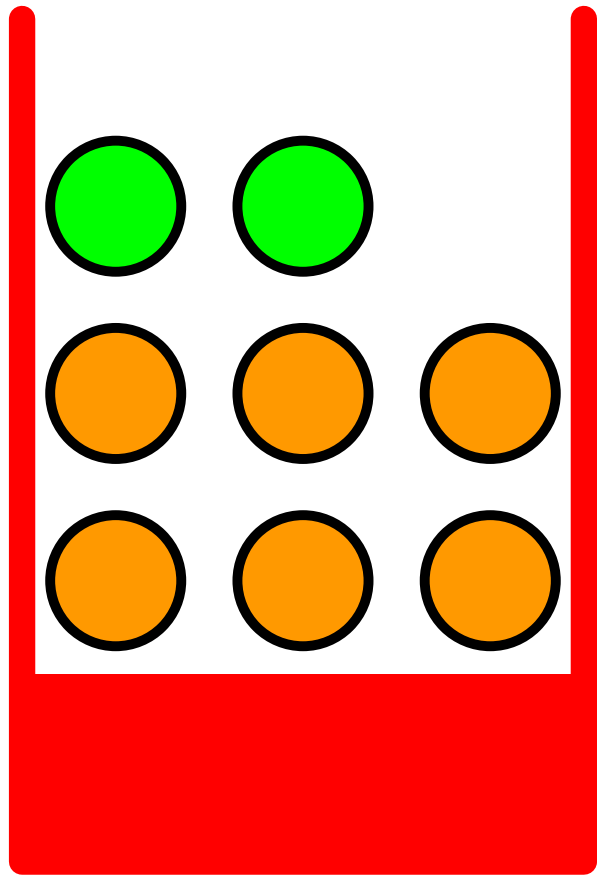
What's the probability of getting a sequence of 1,2,3,4,5,6 if we roll a dice six times?

# Examples: Independent Events

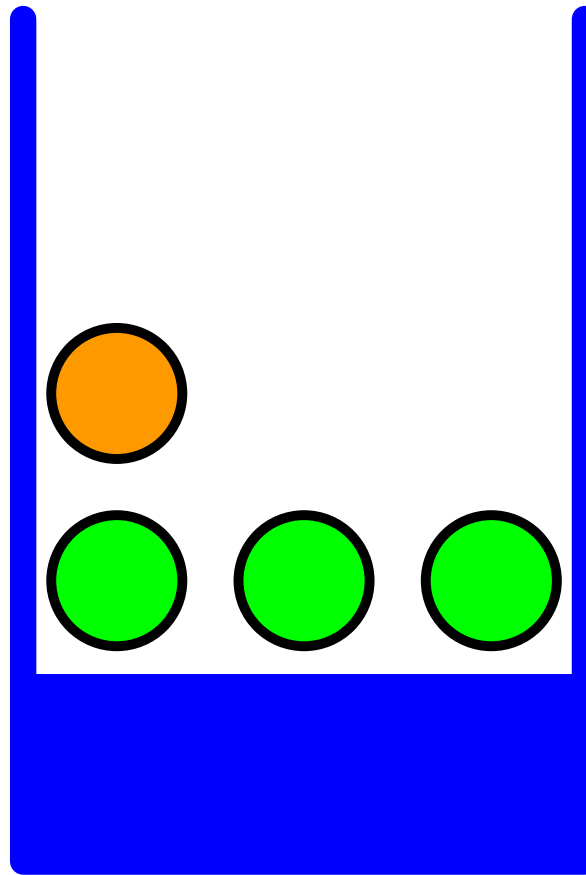
A school survey found that 9 out of 10 students like pizza. If three students are chosen at random with replacement, what is the probability that all three students like pizza?



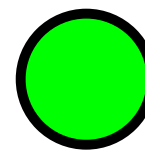
# Urns!



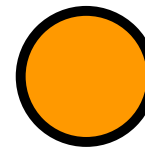
*Red bin*



*Blue bin*

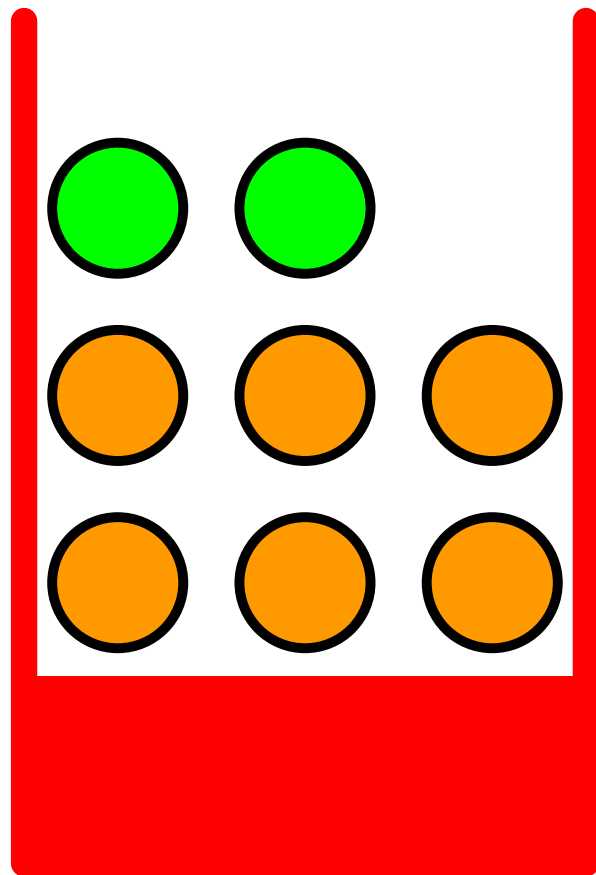


*Apple*

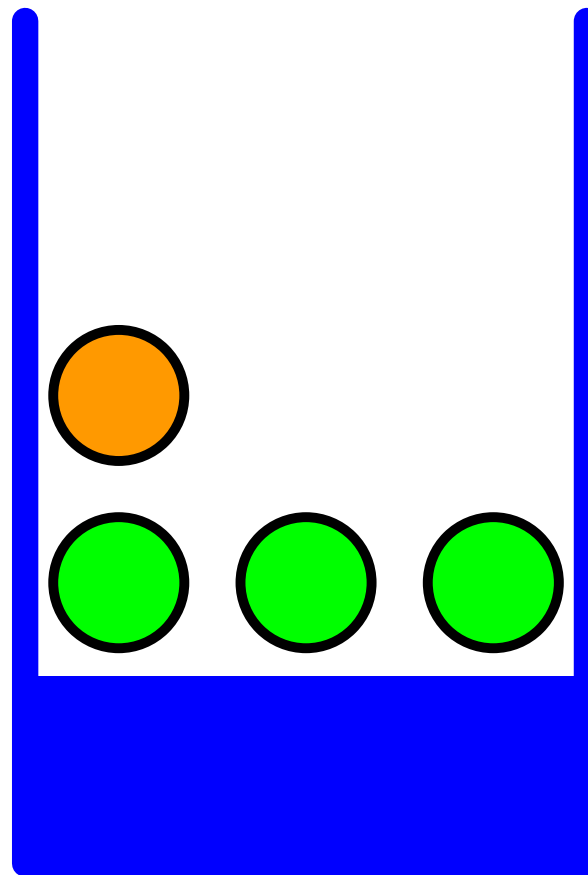


*Orange*

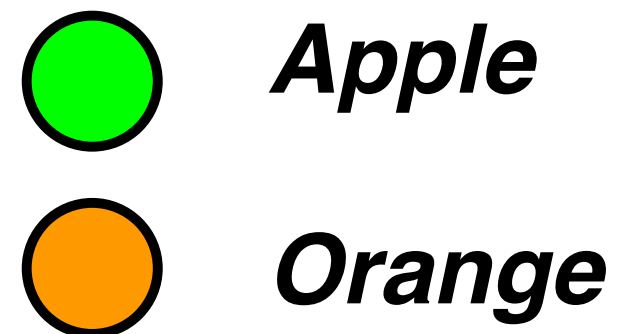
# Dependent Events



***Red bin***

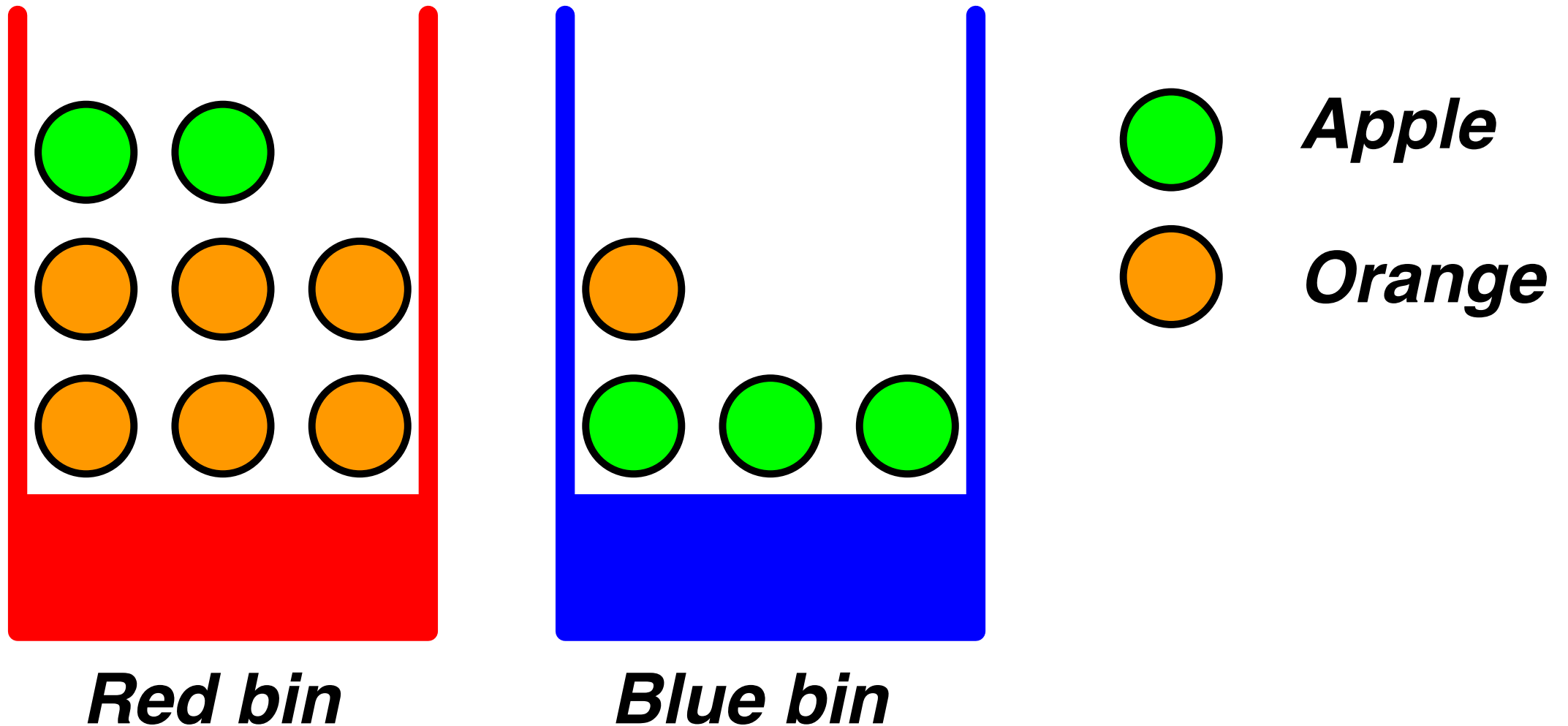


***Blue bin***



*If I randomly pick a fruit from the **red** bin, what is the probability that I get an **apple**?*

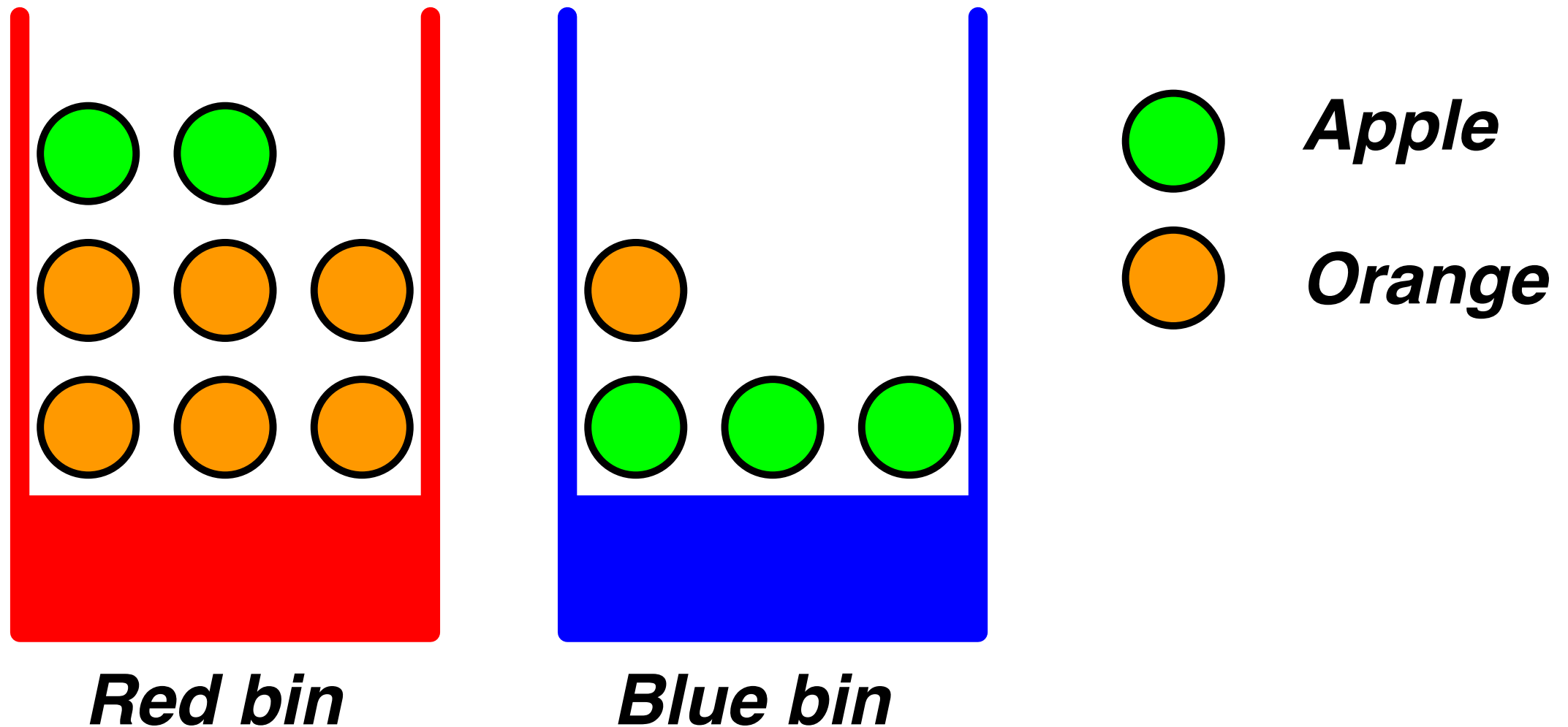
# Dependent Events



*Conditional Probability*

$$P(\text{fruit} = \text{apple} \mid \text{bin} = \text{red}) = 2 / 8$$

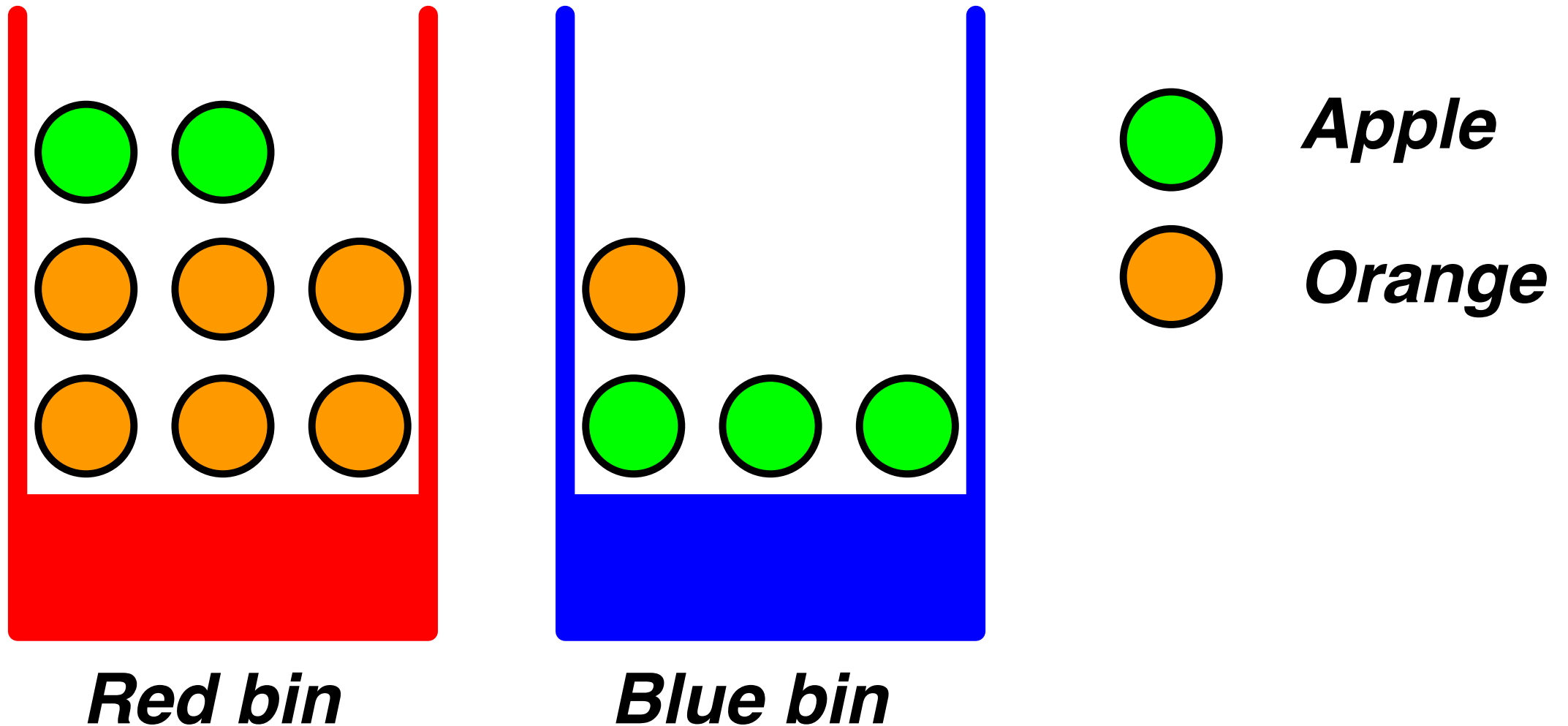
# Dependent Events



*Joint Probability*

$$P(\text{fruit} = \text{apple}, \text{bin} = \text{red}) = 2 / 12$$

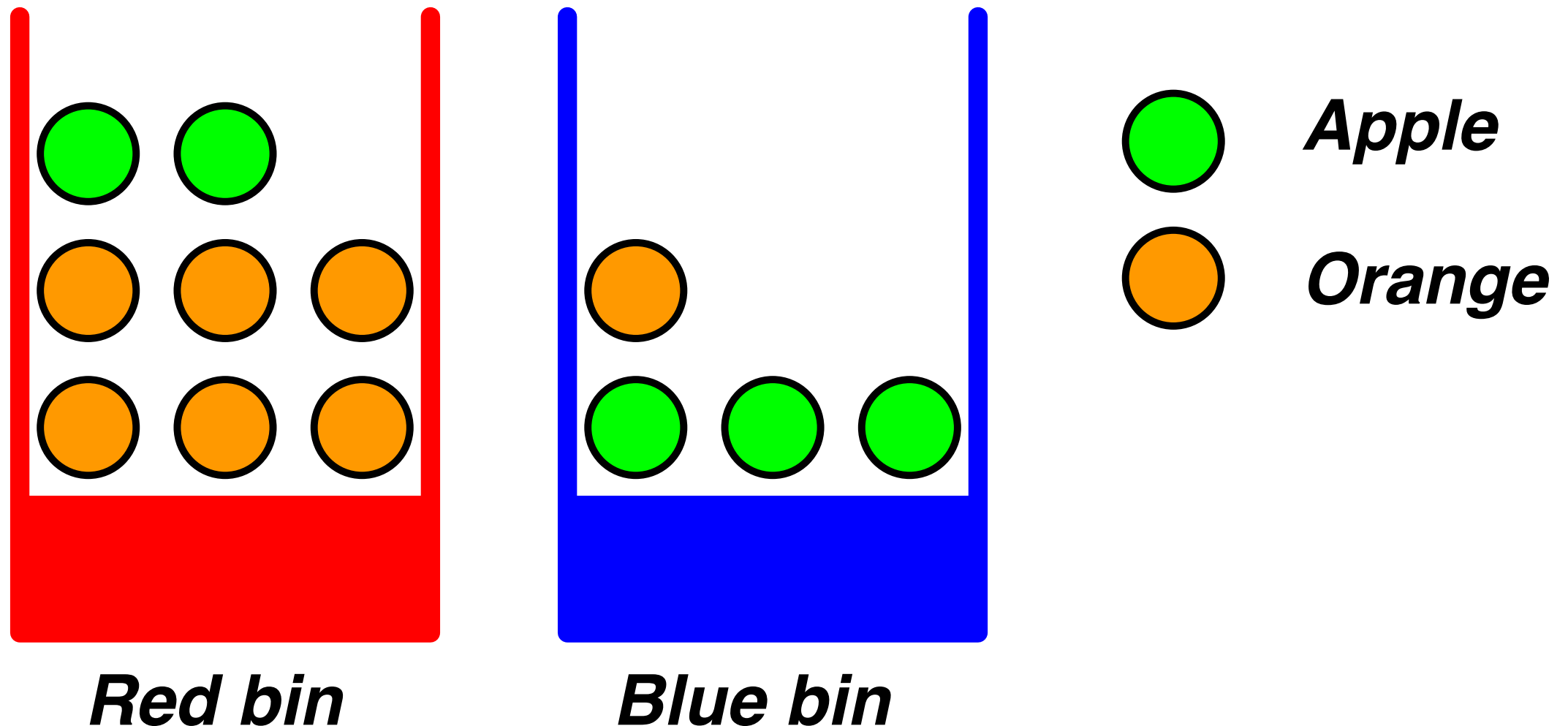
# Dependent Events



*Joint Probability*

$$P(\text{fruit} = \text{apple}, \text{bin} = \text{blue}) = ?$$

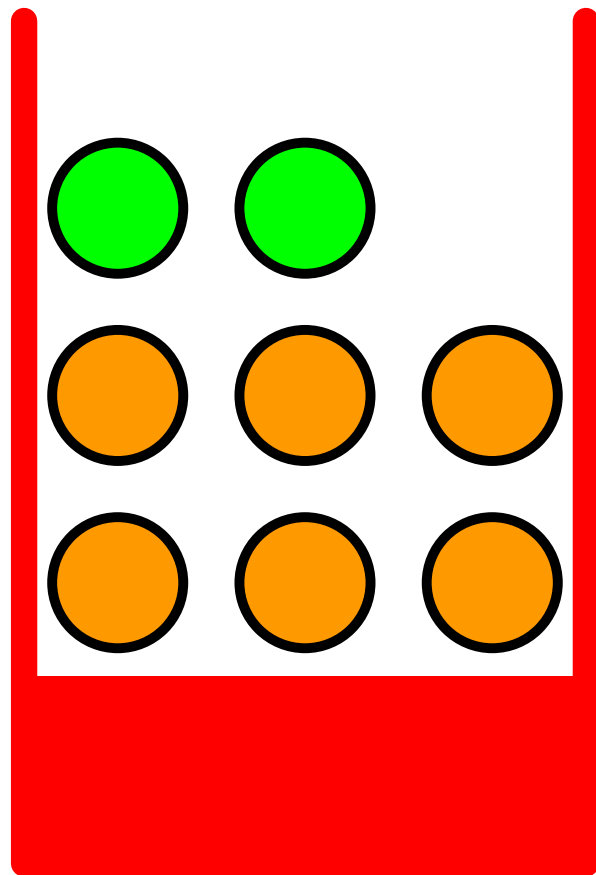
# Dependent Events



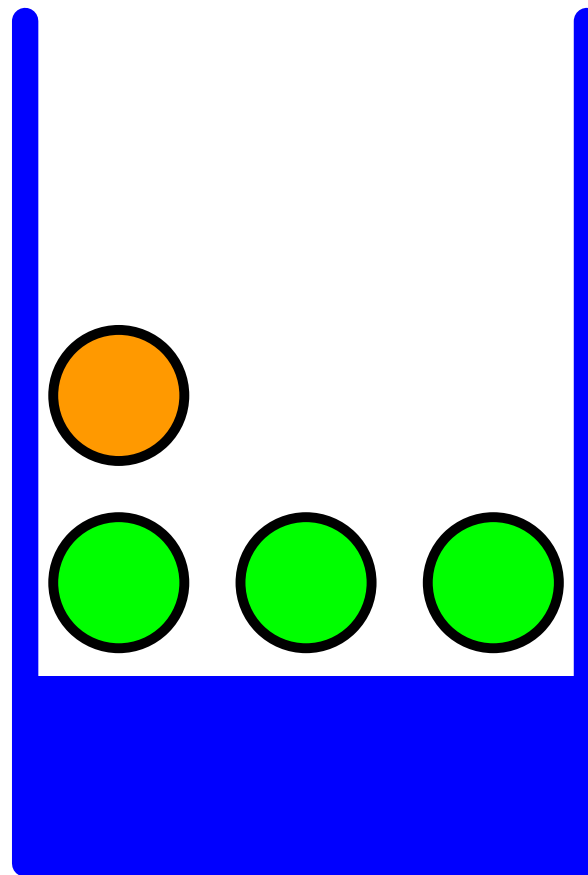
*Joint Probability*

$$P(\text{fruit} = \text{apple}, \text{bin} = \text{blue}) = 3 / 12$$

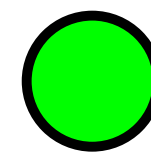
# Dependent Events



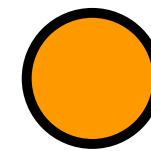
***Red bin***



***Blue bin***



***Apple***

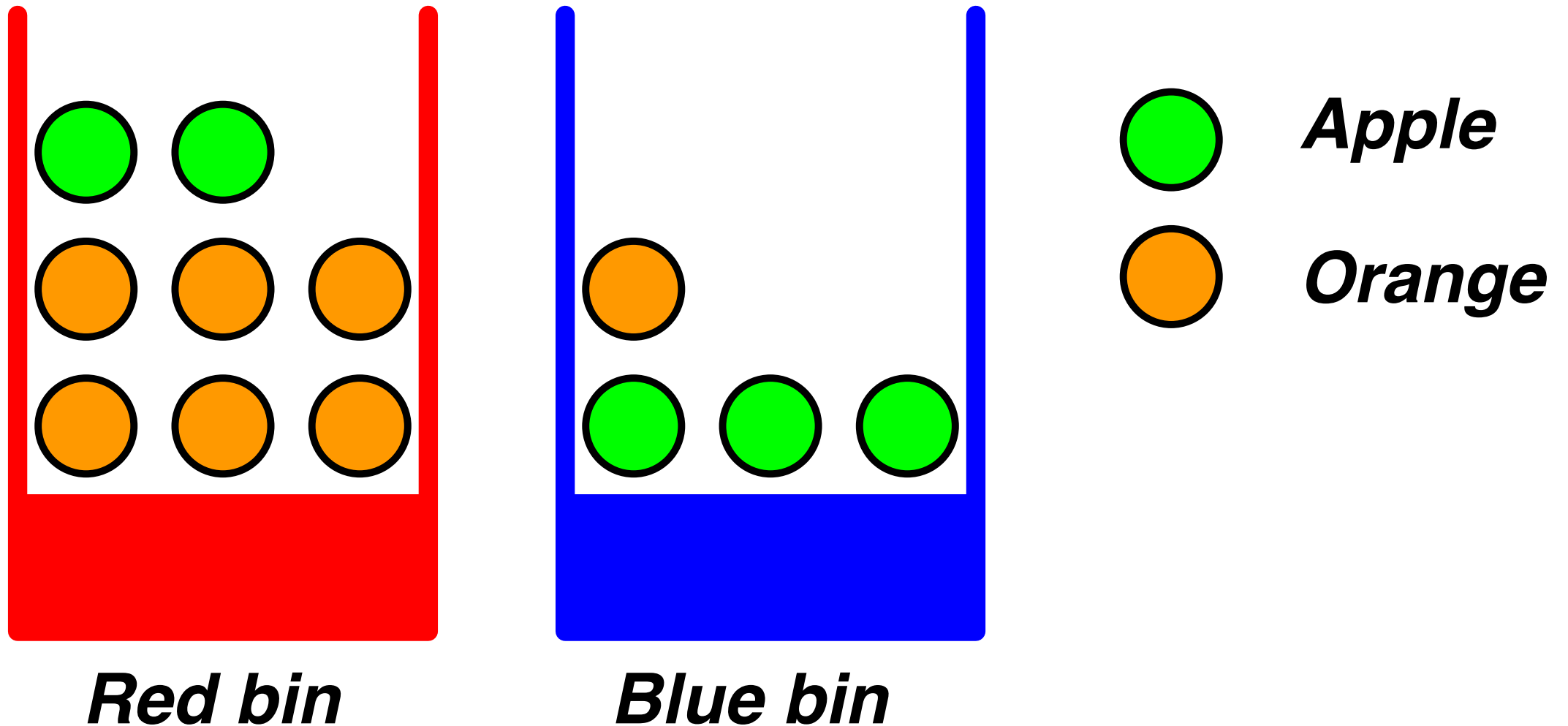


***Orange***

*Joint Probability*

$$P(\text{fruit} = \text{orange}, \text{bin} = \text{blue}) = ?$$

# Dependent Events

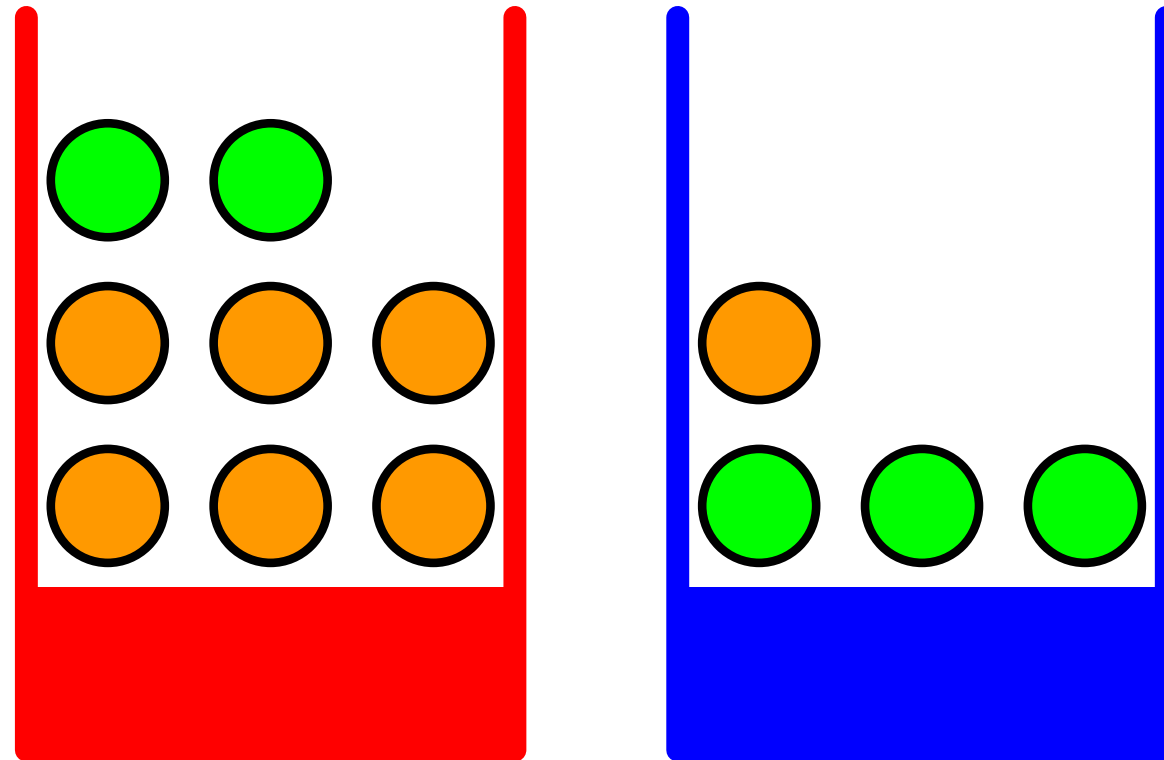


*Joint Probability*

$$P(\text{fruit} = \text{orange}, \text{bin} = \text{blue}) = 1 / 12$$



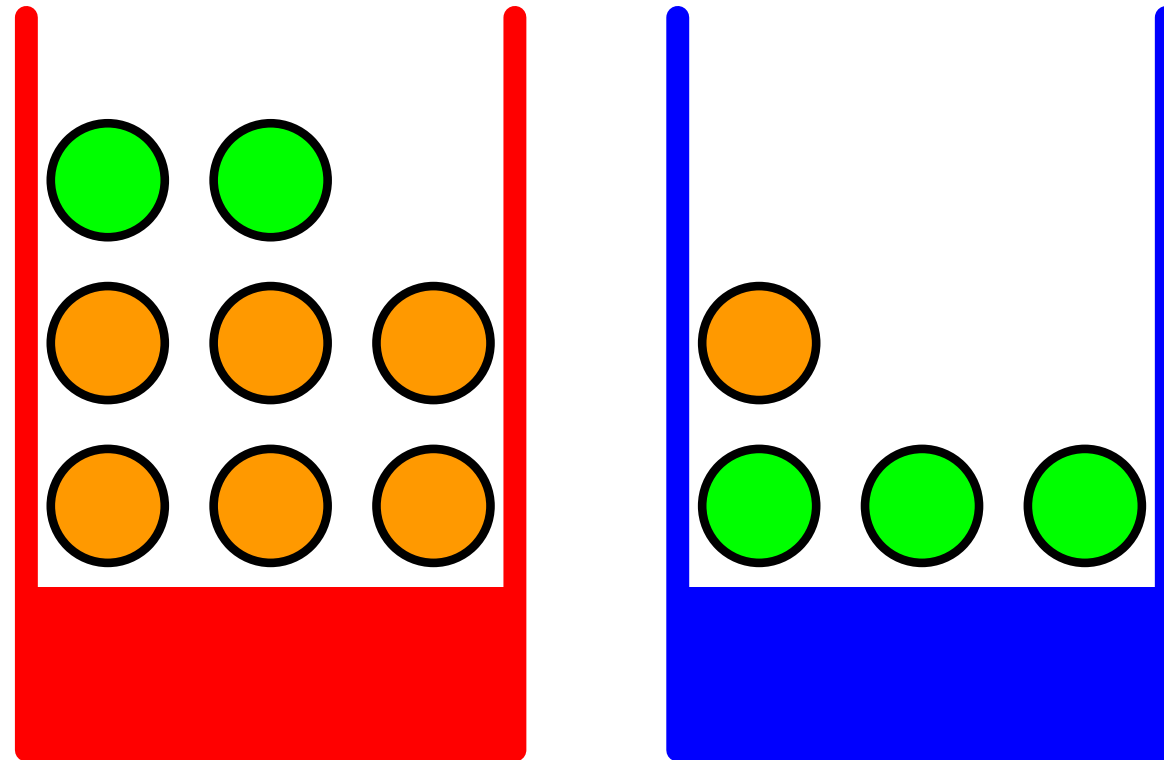
# Two rules of Probability



1. *Sum Rule (Marginal Probabilities)*

$$\begin{aligned} P(\text{fruit} = \text{apple}) &= P(\text{fruit} = \text{apple}, \text{bin} = \text{blue}) \\ &\quad + P(\text{fruit} = \text{apple}, \text{bin} = \text{red}) \\ &= ? \end{aligned}$$

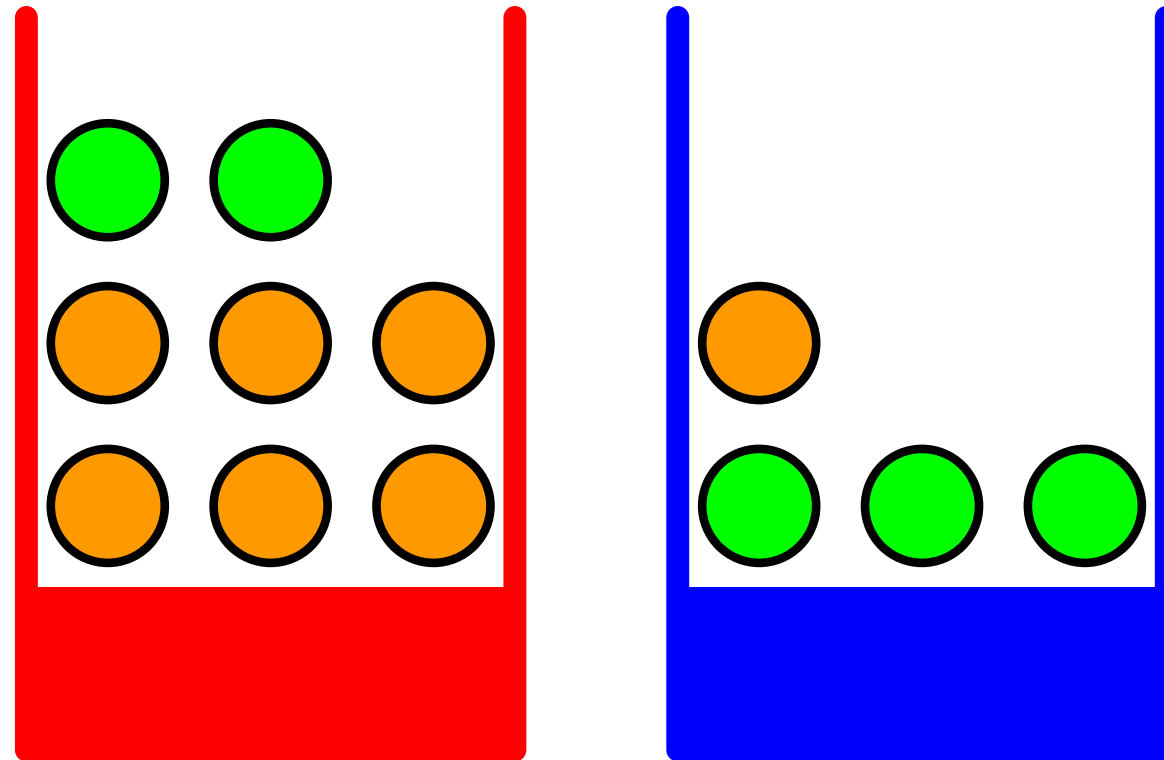
# Two rules of Probability



1. *Sum Rule (Marginal Probabilities)*

$$P(\text{fruit} = \text{apple}) = ?$$

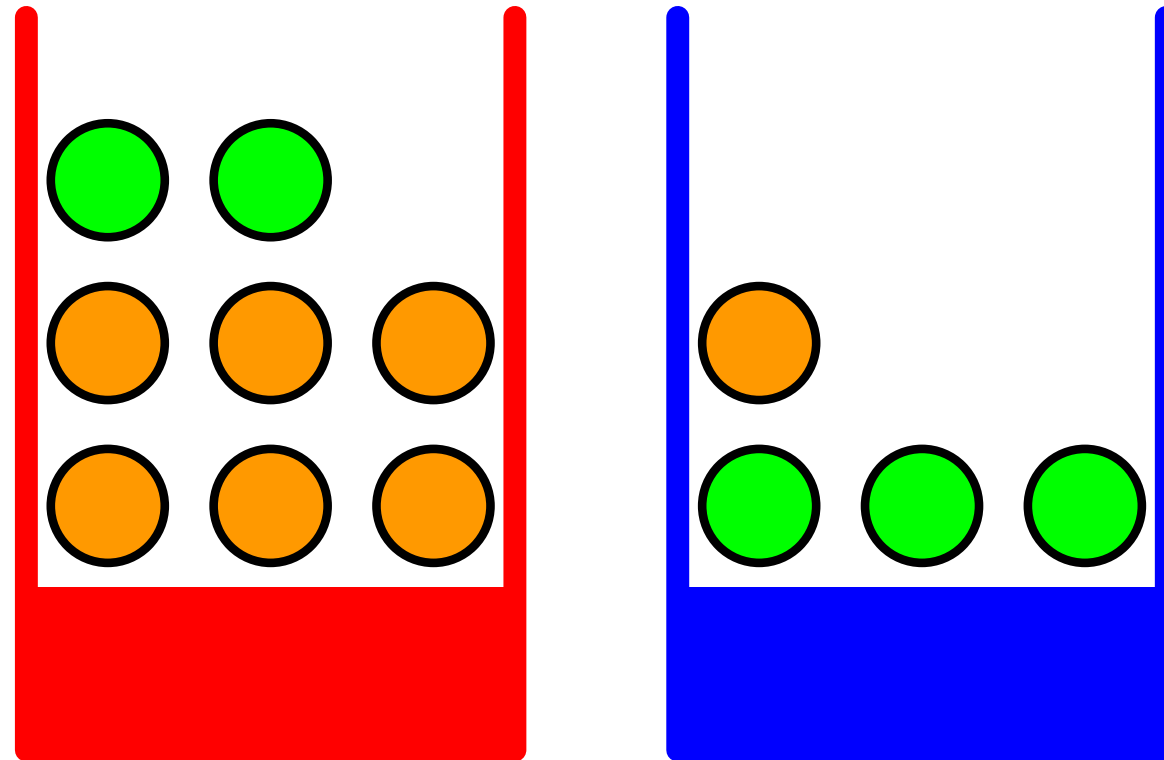
# Two rules of Probability



1. *Sum Rule (Marginal Probabilities)*

$$\begin{aligned} P(\text{fruit} = \text{apple}) &= P(\text{fruit} = \text{apple}, \text{bin} = \text{blue}) \\ &\quad + P(\text{fruit} = \text{apple}, \text{bin} = \text{red}) \\ &= 3 / 12 + 2 / 12 = 5 / 12 \end{aligned}$$

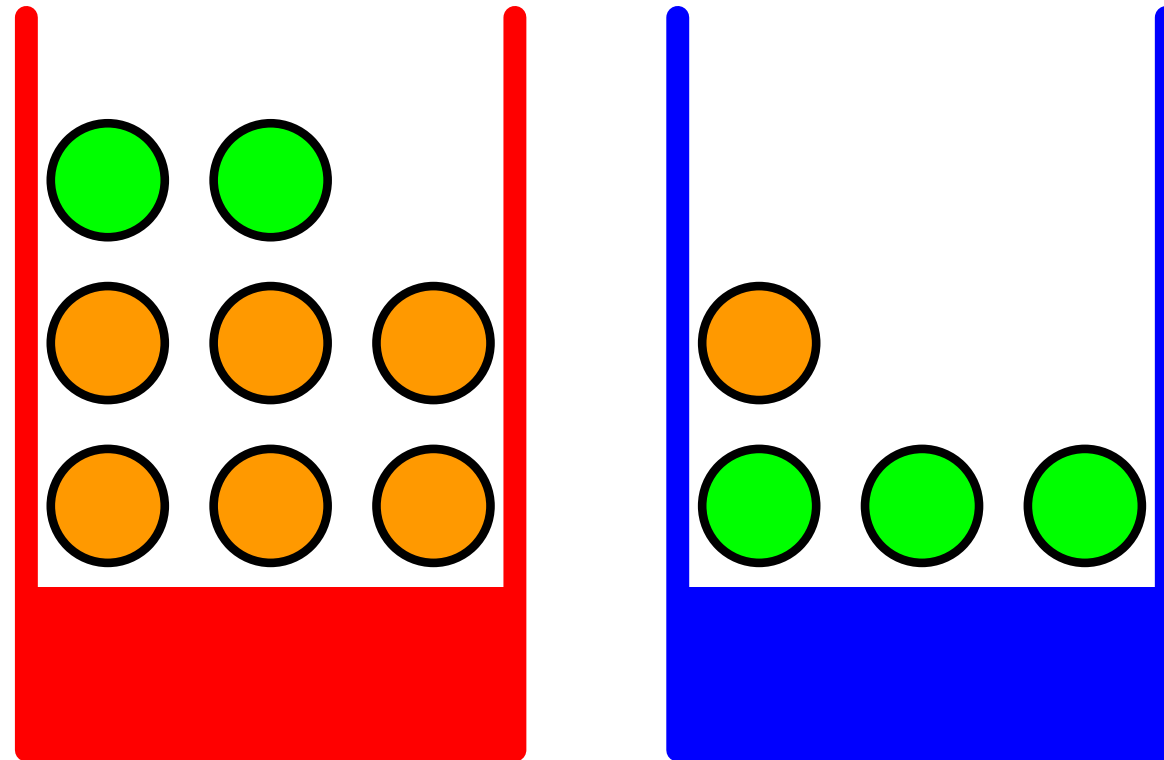
# Two rules of Probability



## *2. Product Rule*

$$P(\text{fruit} = \text{apple}, \text{bin} = \text{red}) = ?$$

# Two rules of Probability



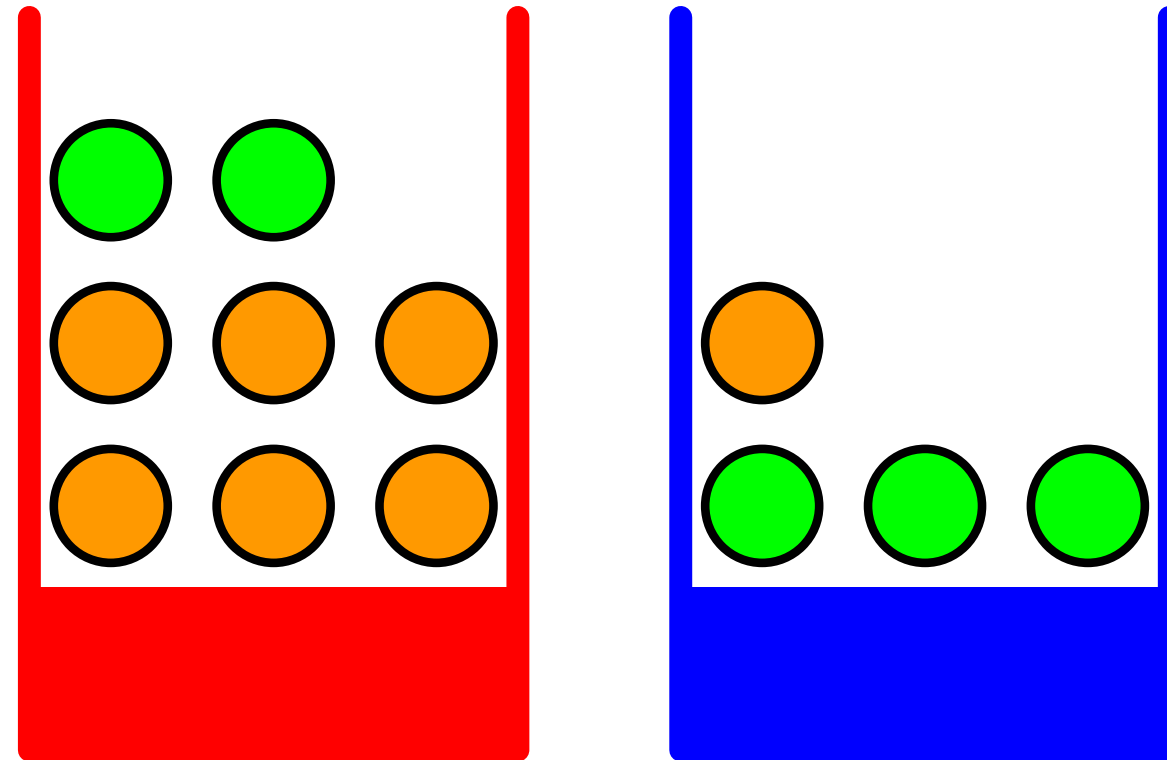
## *2. Product Rule*

$$P(\text{fruit} = \text{apple}, \text{bin} = \text{red}) =$$

$$P(\text{fruit} = \text{apple} \mid \text{bin} = \text{red}) p(\text{bin} = \text{red})$$

$$= ?$$

# Two rules of Probability



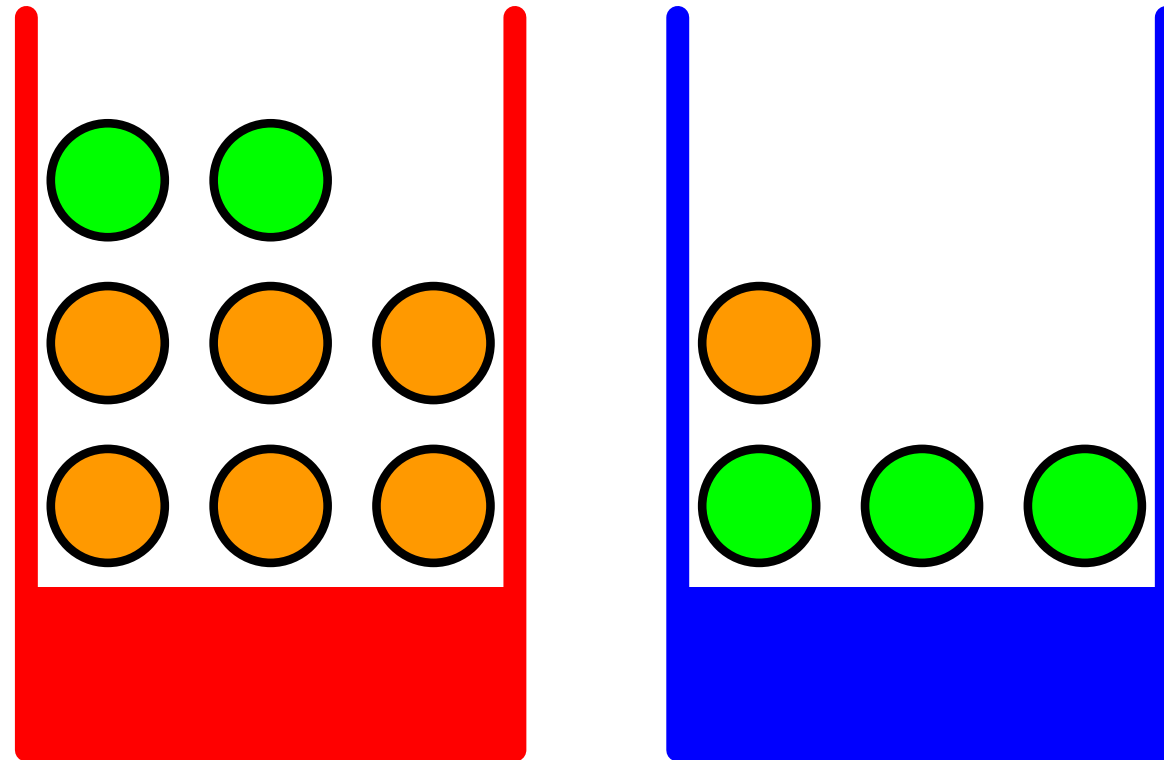
## *2. Product Rule*

$$P(\text{fruit} = \text{apple}, \text{bin} = \text{red}) =$$

$$P(\text{fruit} = \text{apple} \mid \text{bin} = \text{red}) p(\text{bin} = \text{red})$$

$$= 2 / 8 * 8 / 12 = 2 / 12$$

# Two rules of Probability



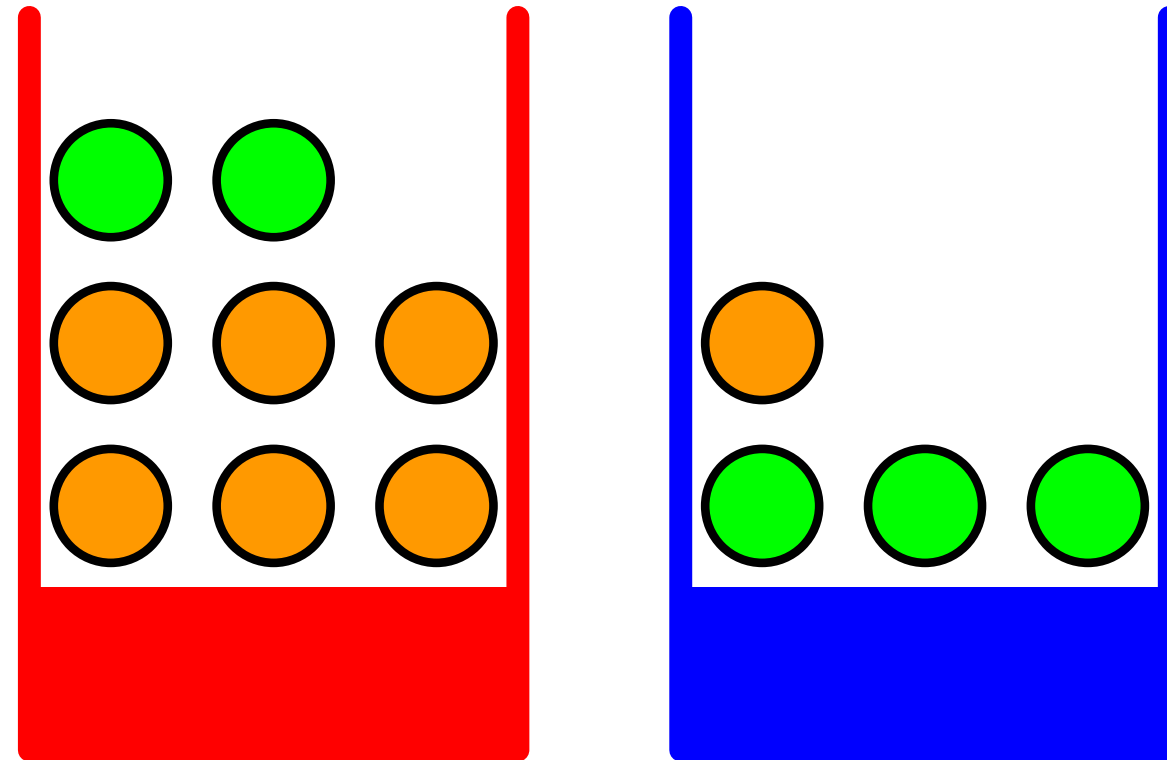
*2. Product Rule (reversed)*

$P(\text{fruit} = \text{apple}, \text{bin} = \text{red}) =$

$P(\text{bin} = \text{red} \mid \text{fruit} = \text{apple}) p(\text{fruit} = \text{apple})$

$= ?$

# Two rules of Probability



*2. Product Rule (reversed)*

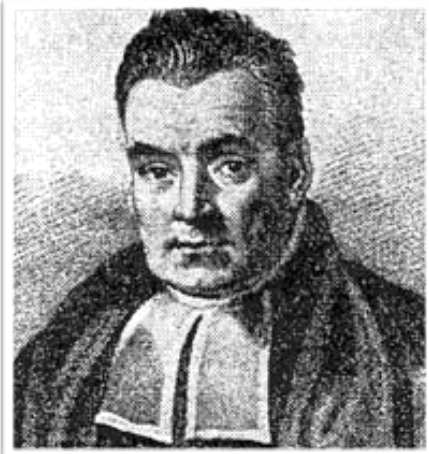
$P(\text{fruit} = \text{apple}, \text{bin} = \text{red}) =$

$P(\text{bin} = \text{red} \mid \text{fruit} = \text{apple}) p(\text{fruit} = \text{apple})$

$= 2 / 5 * 5 / 12 = 2 / 12$



# Bayes' Rule



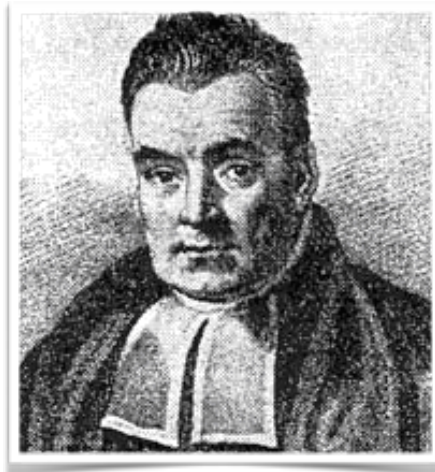
$$p(x | y) = p(y | x)p(x) / p(y)$$

└ Posterior

└ Likelihood

└ Prior

# Bayes' Rule



$$p(\mathbf{x} | \mathbf{y}) = p(\mathbf{y} | \mathbf{x})p(\mathbf{x})/p(\mathbf{y})$$

Posterior

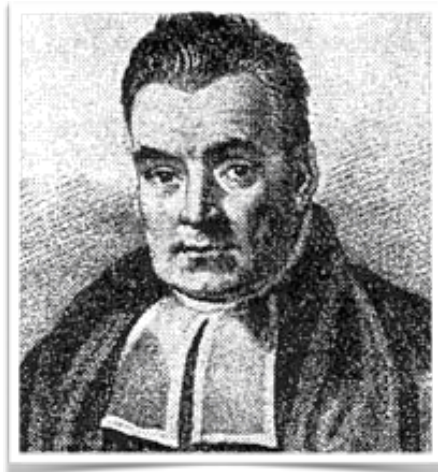
Likelihood

Prior

**Sum Rule:**  $p(\mathbf{y}) = \sum_{\mathbf{x}} p(\mathbf{y}, \mathbf{x})$      $p(\mathbf{x}) = \sum_{\mathbf{y}} p(\mathbf{y}, \mathbf{x})$

**Product Rule:**  $p(\mathbf{y}, \mathbf{x}) = p(\mathbf{y} | \mathbf{x})p(\mathbf{x}) = p(\mathbf{x} | \mathbf{y})p(\mathbf{y})$

# Bayes' Rule



$$p(\mathbf{x} | \mathbf{y}) = p(\mathbf{y} | \mathbf{x})p(\mathbf{x})/p(\mathbf{y})$$

Posterior

Likelihood

Prior

$p(\mathbf{x})$

*Probability of rare disease: 0.005*

$p(\mathbf{y} | \mathbf{x})$

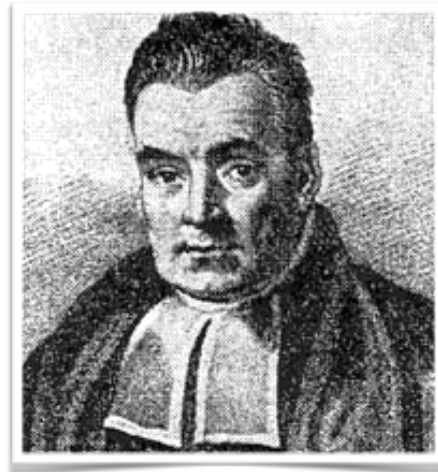
*Probability of detection: 0.98*

*Probability of false positive: 0.05*

$p(\mathbf{x} | \mathbf{y})$

*Probability of disease when test positive?*

# Bayes' Rule



$$p(\mathbf{x} | \mathbf{y}) = p(\mathbf{y} | \mathbf{x})p(\mathbf{x})/p(\mathbf{y})$$

└ Posterior

└ Likelihood

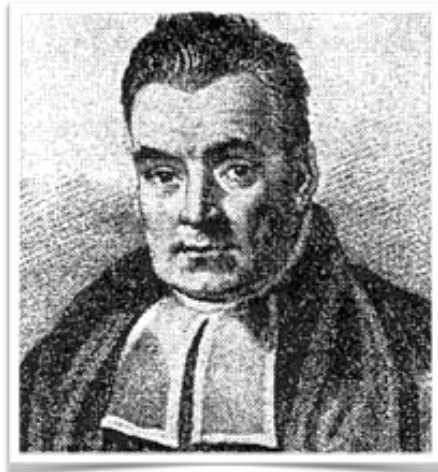
└ Prior

$$p(\mathbf{y}, \mathbf{x}) = p(\mathbf{y} | \mathbf{x})p(\mathbf{x})$$

$$p(\mathbf{y}) = \sum_{\mathbf{x}} p(\mathbf{y}, \mathbf{x})$$

$$p(\mathbf{x} | \mathbf{y})$$

# Bayes' Rule



$$p(\mathbf{x} | \mathbf{y}) = p(\mathbf{y} | \mathbf{x})p(\mathbf{x}) / p(\mathbf{y})$$

Posterior

Likelihood

Prior

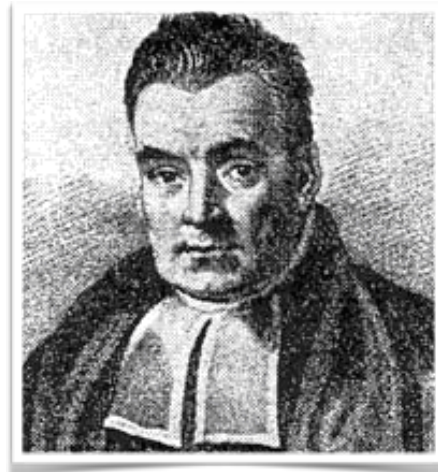
$$p(\mathbf{y}, \mathbf{x}) = p(\mathbf{y} | \mathbf{x})p(\mathbf{x})$$

$$0.98 * 0.005 = 0.0049$$

$$p(\mathbf{y}) = \sum_{\mathbf{x}} p(\mathbf{y}, \mathbf{x})$$

$$p(\mathbf{x} | \mathbf{y})$$

# Bayes' Rule



$$p(\mathbf{x} | \mathbf{y}) = p(\mathbf{y} | \mathbf{x})p(\mathbf{x})/p(\mathbf{y})$$

Posterior

Likelihood

Prior

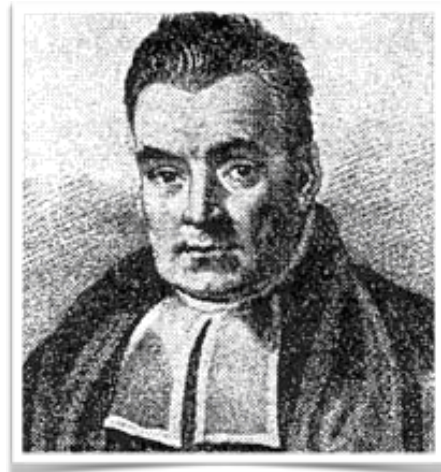
$$p(\mathbf{y}, \mathbf{x}) = p(\mathbf{y} | \mathbf{x})p(\mathbf{x})$$

$$0.98 * 0.0005 = 0.00049$$

$$p(\mathbf{y}) = \sum_{\mathbf{x}} p(\mathbf{y}, \mathbf{x}) \quad 0.98 * 0.0005 + 0.05 * 0.9995 = 0.05447$$

$$p(\mathbf{x} | \mathbf{y})$$

# Bayes' Rule



$$p(\mathbf{x} | \mathbf{y}) = p(\mathbf{y} | \mathbf{x})p(\mathbf{x}) / p(\mathbf{y})$$

Posterior

Likelihood

Prior

$$p(\mathbf{y}, \mathbf{x}) = p(\mathbf{y} | \mathbf{x})p(\mathbf{x})$$

$$0.98 * 0.0005 = 0.00049$$

$$p(\mathbf{y}) = \sum_{\mathbf{x}} p(\mathbf{y}, \mathbf{x}) \quad 0.98 * 0.0005 + 0.05 * 0.9995 = 0.05447$$

$$p(\mathbf{x} | \mathbf{y})$$

$$0.00049 / 0.05447 = 0.00899$$

# Random Variables

- ***Random Variable:*** A variable with a stochastic *outcome*

$$X = x \quad x \in \{1, 2, 3, 4, 5, 6\}$$



# Random Variables

- **Random Variable:** A variable with a stochastic *outcome*

$$X = x \quad x \in \{1, 2, 3, 4, 5, 6\}$$

- **Event:** A set of *outcomes*

$$X \geq 3 \quad \{3, 4, 5, 6\}$$

# Random Variables

- **Random Variable:** A variable with a stochastic *outcome*

$$X = x \quad x \in \{1, 2, 3, 4, 5, 6\}$$

- **Event:** A set of *outcomes*

$$X \geq 3 \quad \{3, 4, 5, 6\}$$

- **Probability:** The chance that a randomly selected *outcome* is part of an *event*

$$P(X \geq 3) = 4 / 6$$

# Distribution

- A *distribution* maps *outcomes* to *probabilities*

$$P(X = x) = 1 / 6$$

- Commonly used (or abused) shorthand:

$$P(x) \text{ is equivalent to } P(X = x)$$

# Probability Spaces

**Definition:** A probability space  $(\Omega, \mathcal{F}, \mathbb{P})$  consists of

- A sample space  $\Omega$  (i.e. the set of *outcomes*)
- A set of events  $\mathcal{F}$  (i.e. the set possible sets)
- A probability measure  $\mathbb{P}$  (maps events to probabilities)

# Probability Spaces

**Definition:** A probability space  $(\Omega, \mathcal{F}, P)$  consists of

- A sample space  $\Omega$  (i.e. the set of *outcomes*)
- A set of events  $\mathcal{F}$  (i.e. the set possible sets)
- A probability measure  $P$  (maps events to probabilities)

## Axioms of Probability

$$P : \mathcal{F} \rightarrow \mathbb{R} \quad P(E) \geq 0 \quad \forall E \in \mathcal{F} \quad P(\Omega) = 1$$

$$P(E_1, E_2) = P(E_1) + P(E_2) \quad \text{when } E_1 \cap E_2 = \emptyset$$

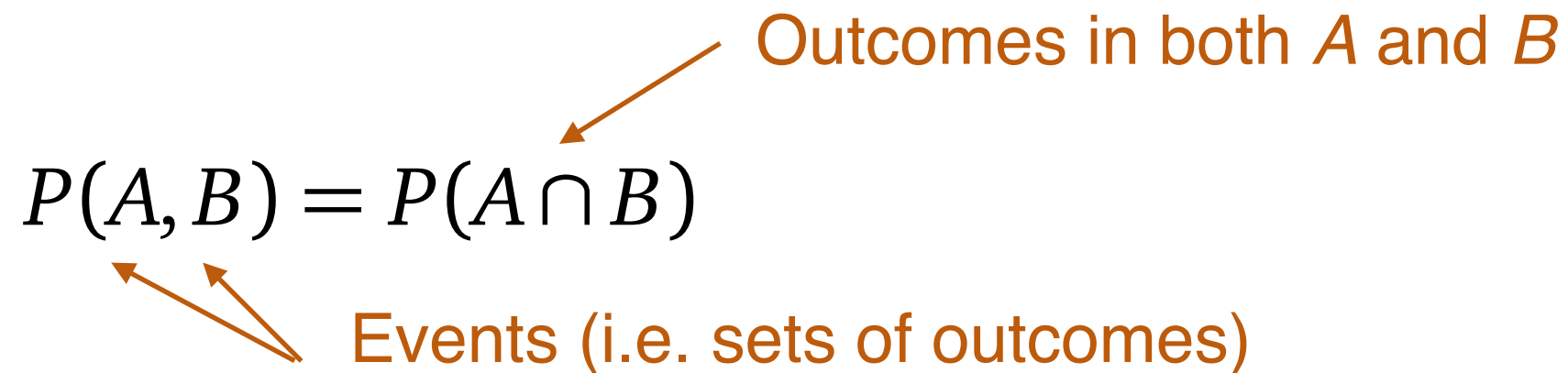
# Conditional Probabilities

- **Definition:** Joint Probability

$$P(A, B) = P(A \cap B)$$

Outcomes in both  $A$  and  $B$

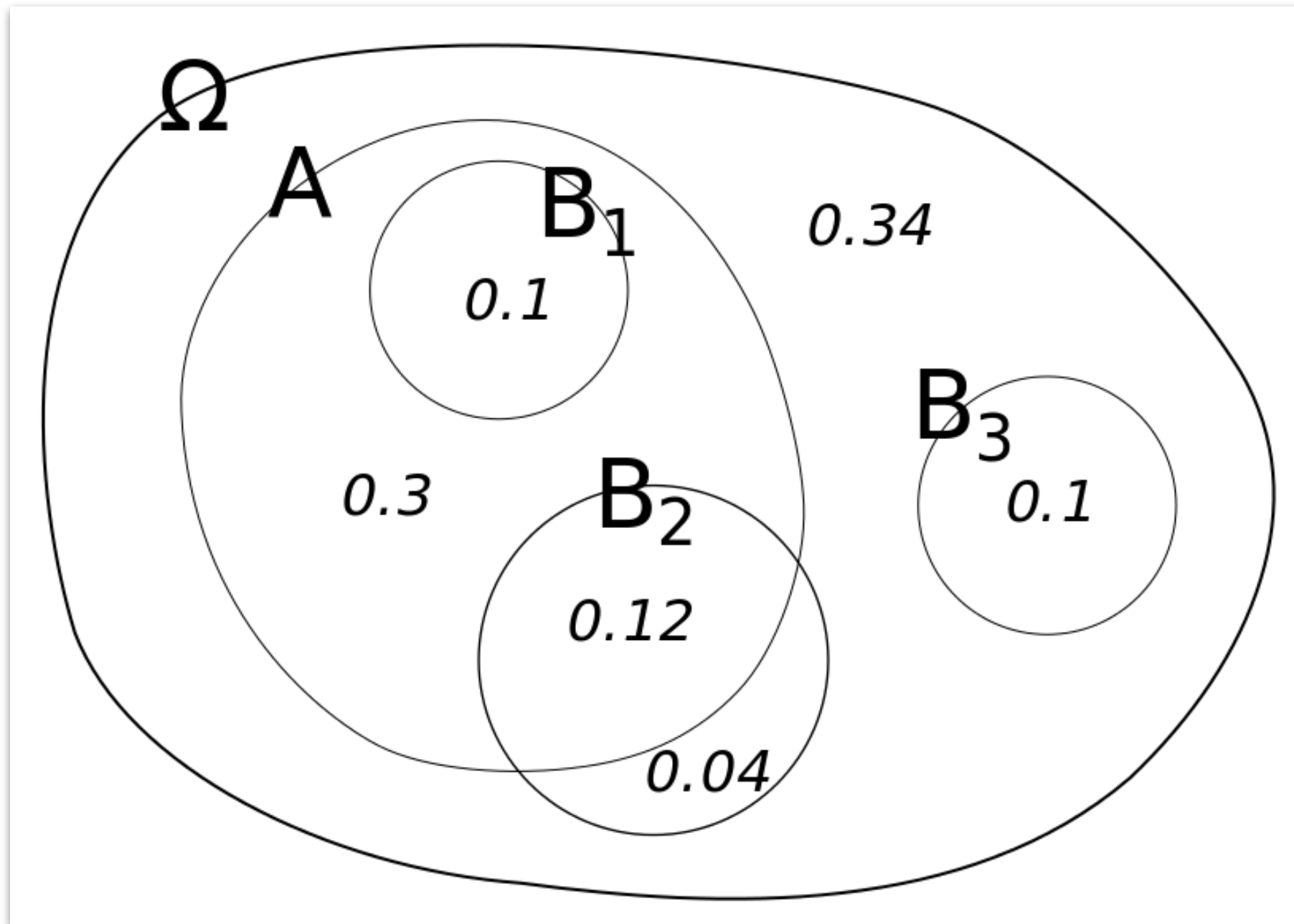
Events (i.e. sets of outcomes)



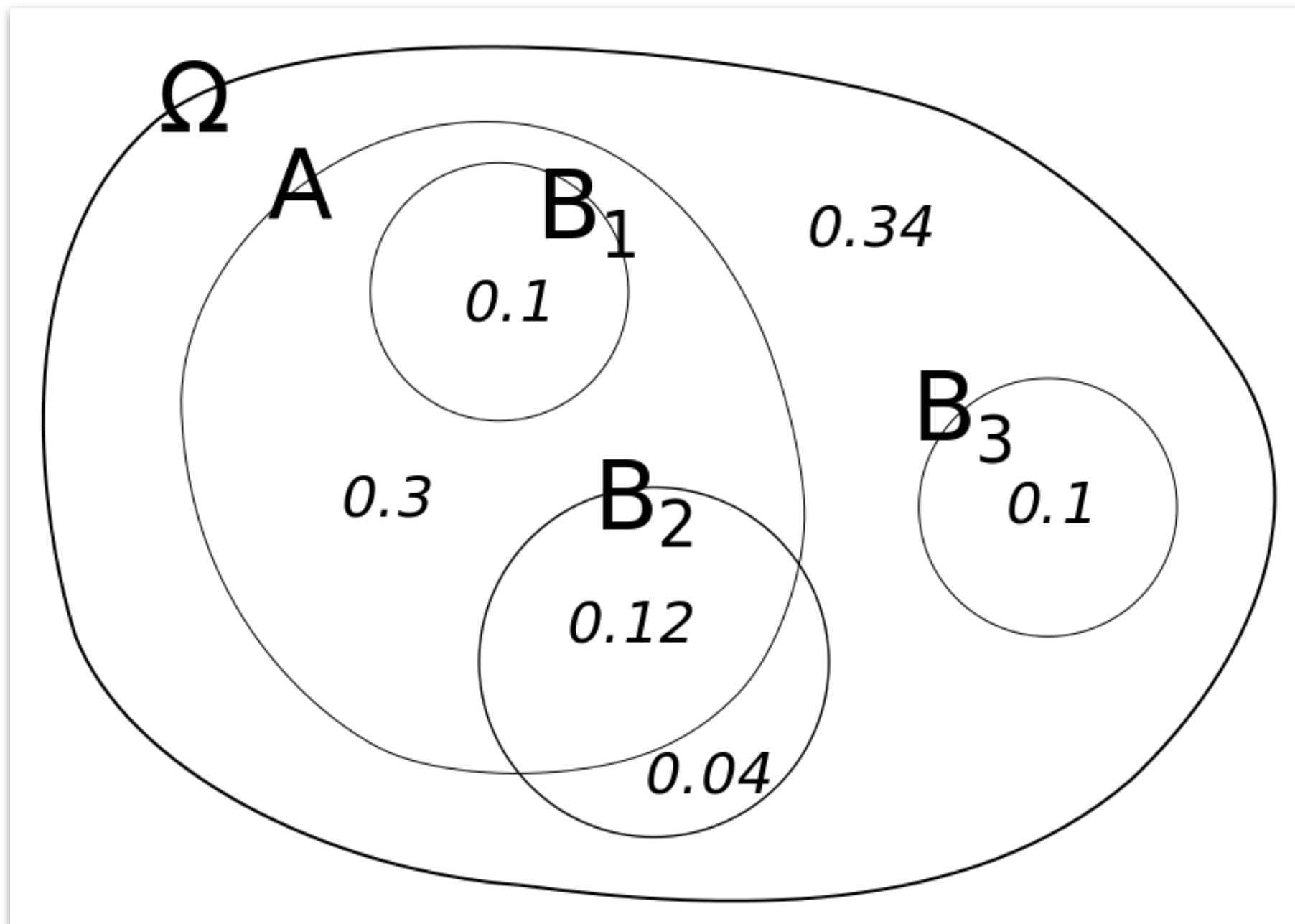
- **Definition:** Conditional Probability

$$P(A | B) = \frac{P(A, B)}{P(B)}$$

# Conditional Probability



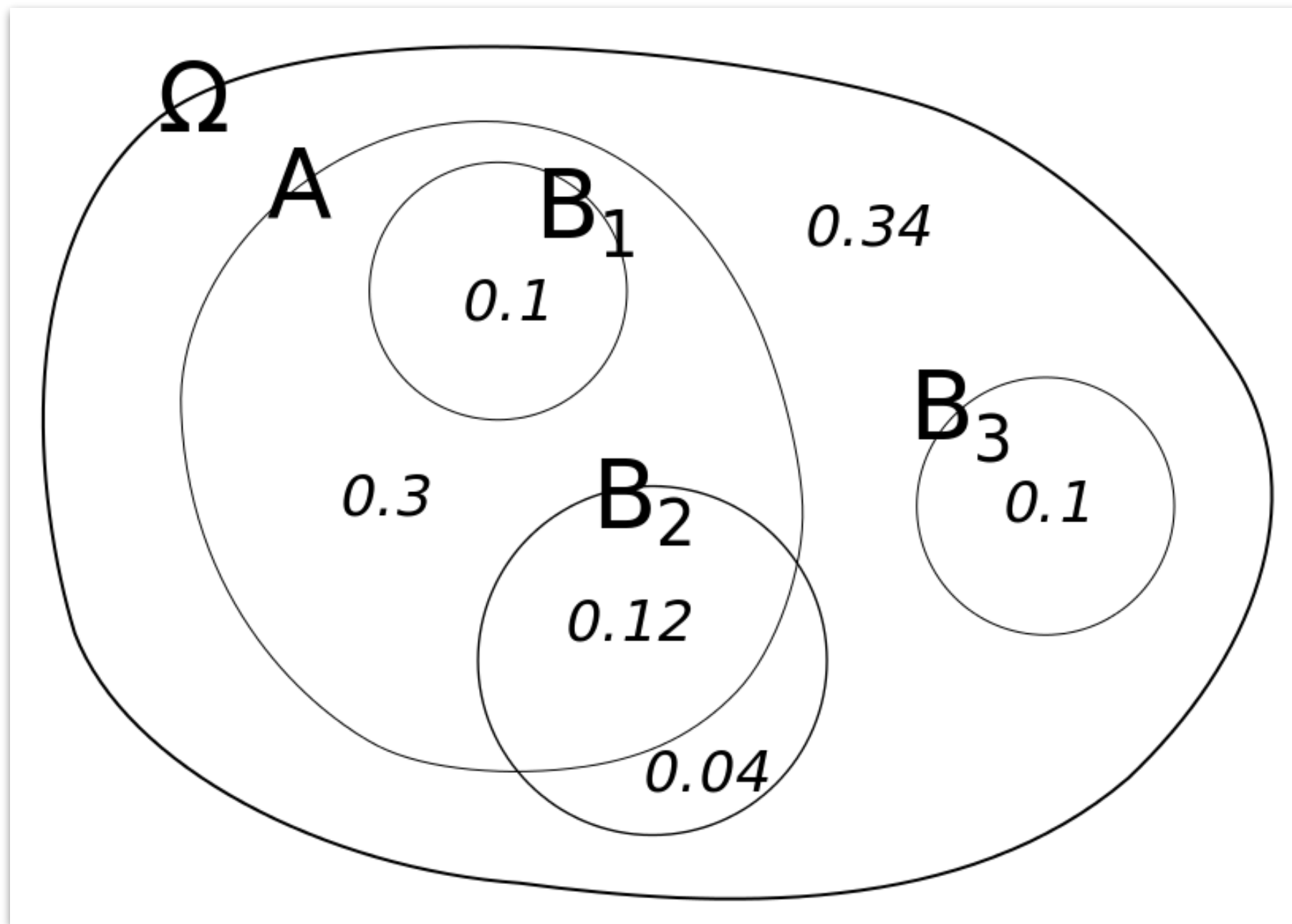
# Conditional Probability



What is the probability  $P(B_3)$ ? **0.1 / 0.34**

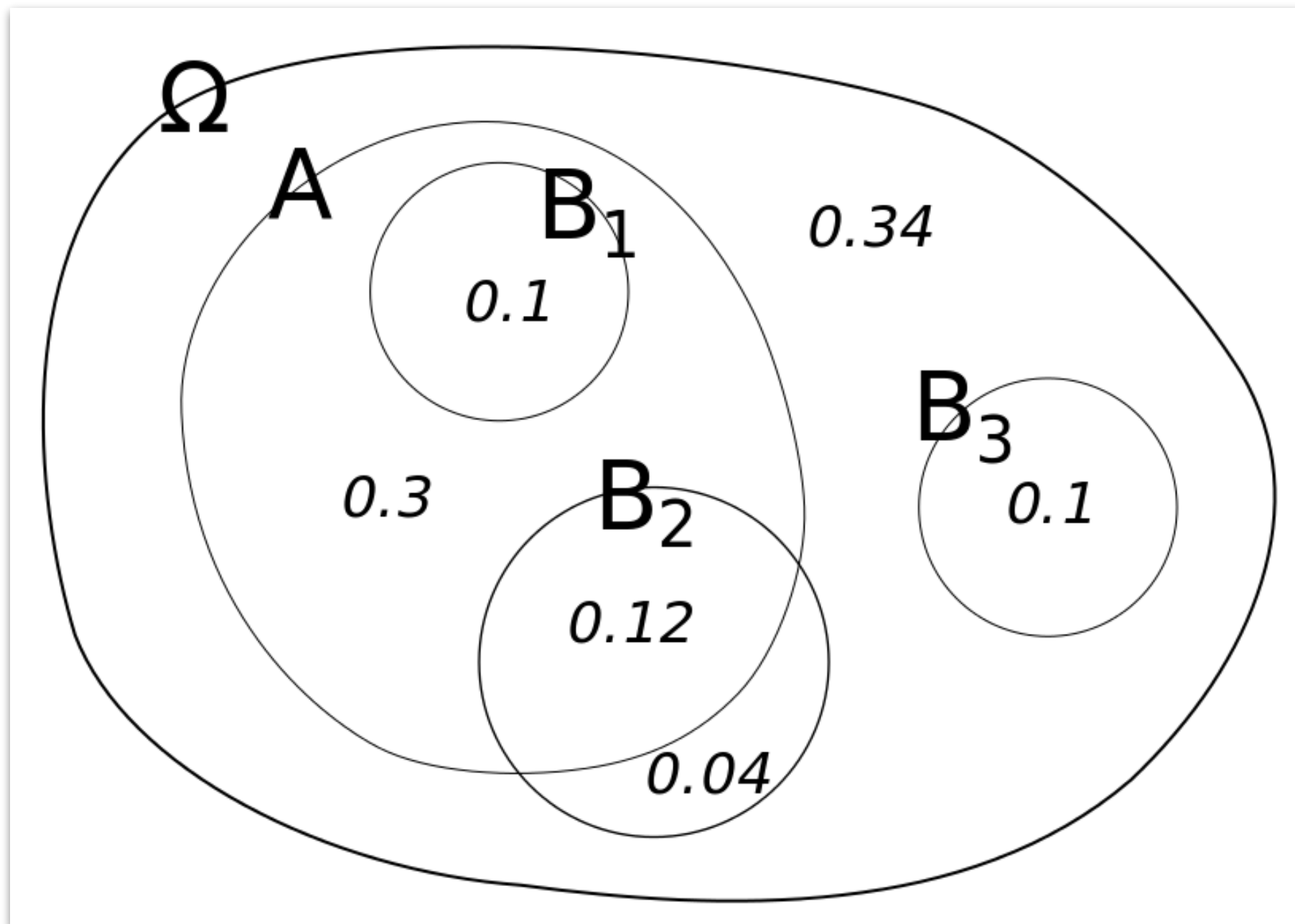


# Conditional Probability



What is the probability  $P(B_2 \mid A)$ ?  **$0.12 / 0.3$**

# Conditional Probability



What is the probability  $P(B_1 \mid B_3)$ ? **0.0 / 0.1**

# Examples: Conditional Probability

1. A math teacher gave her class two tests.
  - *25% of the class passed both tests*
  - *42% of the class passed the first test.*

*What percent of those who passed the first test also passed the second test?*

# Examples: Conditional Probability

2. Suppose that for houses in New England
- *84% of the houses have a garage*
  - *65% of the houses have a garage and a backyard.*

*What is the probability that a house has a backyard given that it has a garage?*

To Jupyter...

# Probability Density Functions

- **Problem:** If  $X$  is a *continuous* variable, then  $P(X=x)$  is 0 for any outcome  $x$

$$X \sim \text{Normal}(0, 1)$$

$$P(X = \pi) = 0$$

$$P(3.1 \leq X \leq 3.2) \neq 0$$

Single Outcome

Event

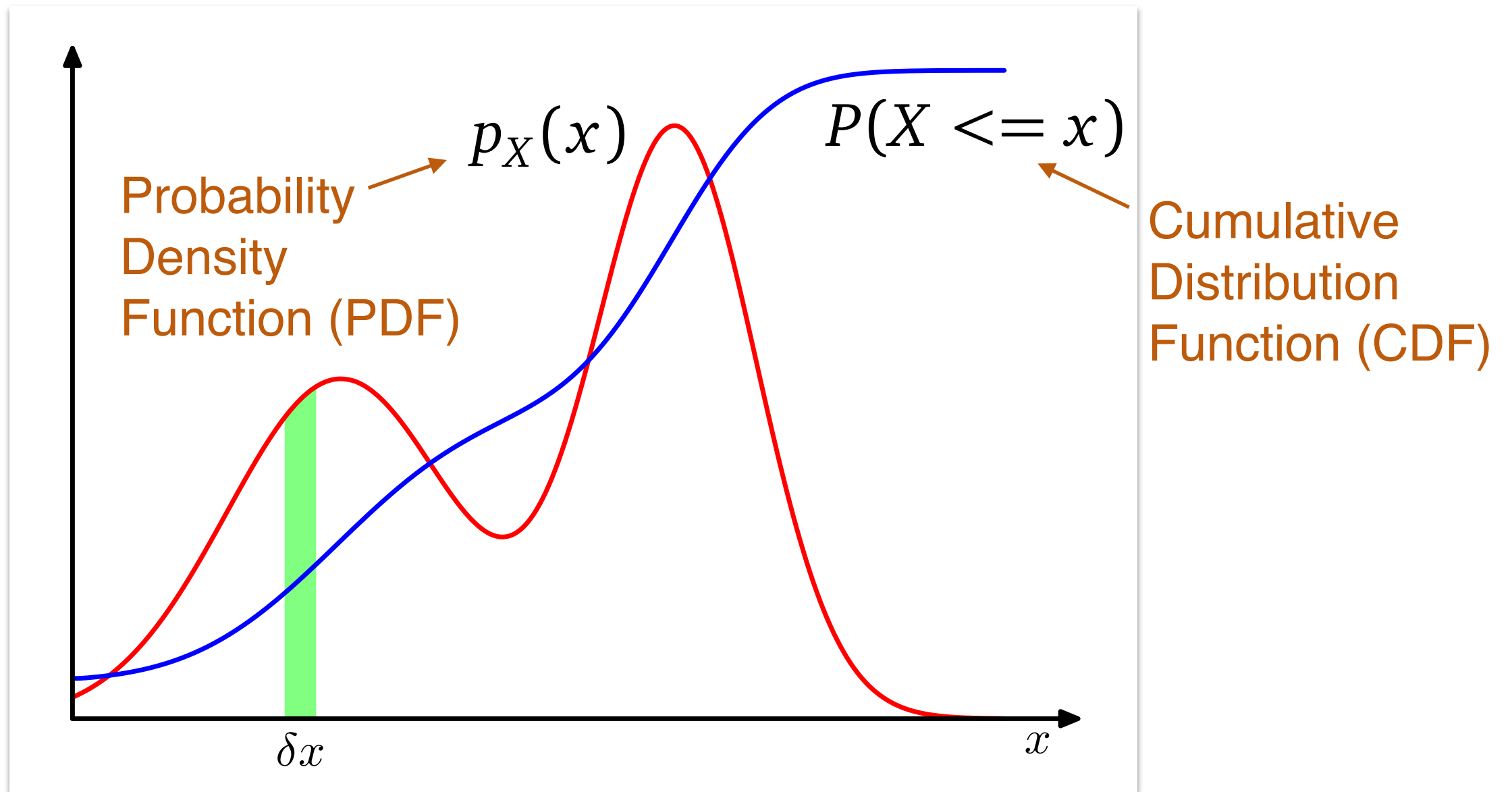
- **Solution:** Define a density function as a derivative

$$p_X(x) = \lim_{\delta \rightarrow 0} \frac{P(x - \delta < X < x + \delta)}{2\delta}$$

Capital P for probability

Small p for density

# Probability Density Functions



$$p_X(x) = \lim_{\delta \rightarrow 0} \frac{P(x - \delta < X < x + \delta)}{2\delta}$$

# Expected Values

$X \sim p(x)$  ←  $X$  is a random variable with density  $p(x)$

## Statistics

$$\mathbb{E}[X] = \sum_x p(x) x$$

$$\mathbb{E}[X] = \int dx p(x) x$$

## Machine Learning

$$\mathbb{E}_{p(x|y)}[f(x)] = \sum_x p(x|y) f(x)$$

$$\mathbb{E}_{p(x|y)}[f(x)] = \int dx p(x|y) f(x)$$



# Mean, Variance, Covariance

## Mean

$$\mu_X = \mathbb{E}[X]$$

## Variance

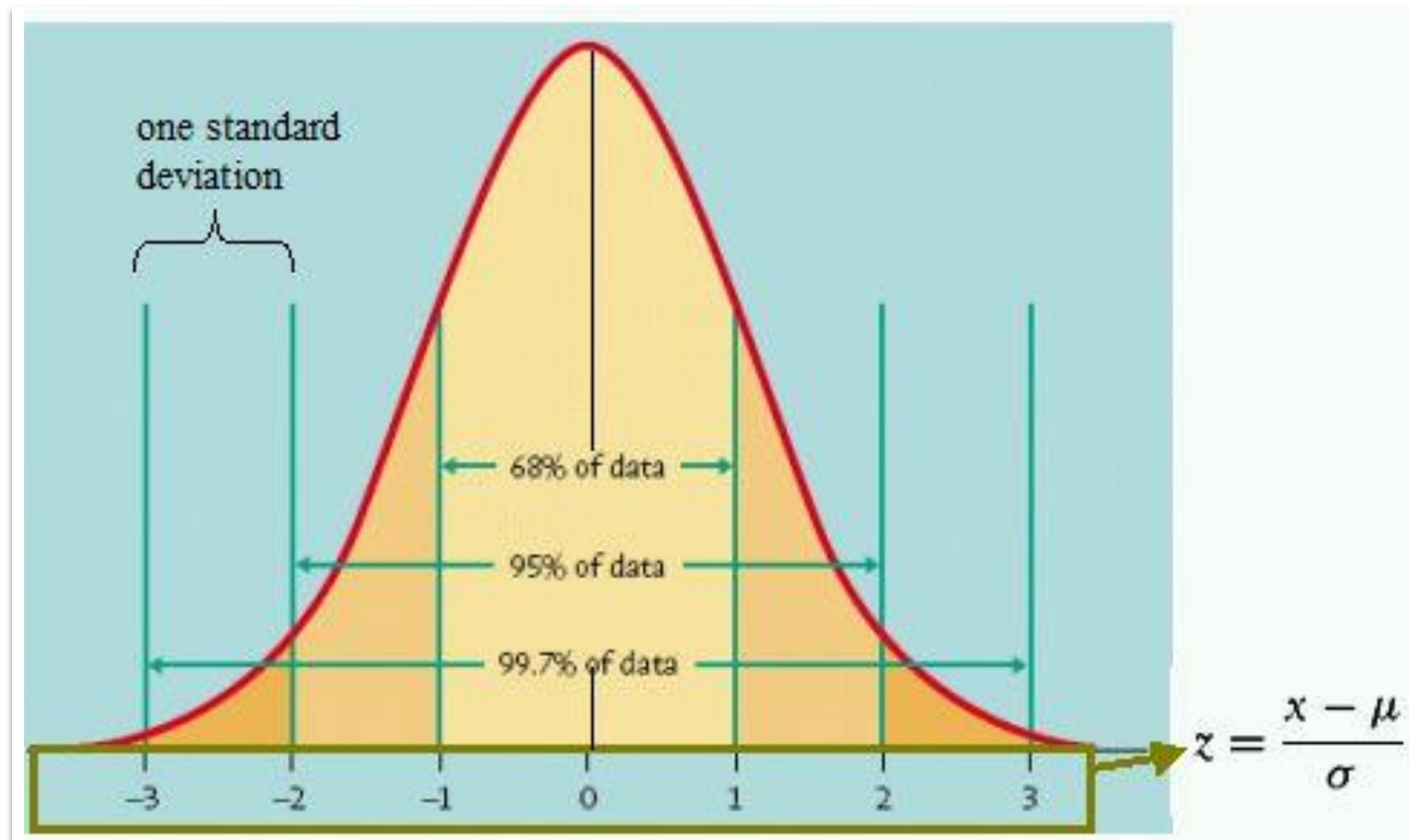
$$\sigma_X^2 = \text{Var}[X] = \mathbb{E}[(X - \mu_X)^2]$$

## Covariance

$$\Sigma_{X,Y} = \text{Cov}[X, Y] = \mathbb{E}[(X - \mu_X)(Y - \mu_Y)]$$

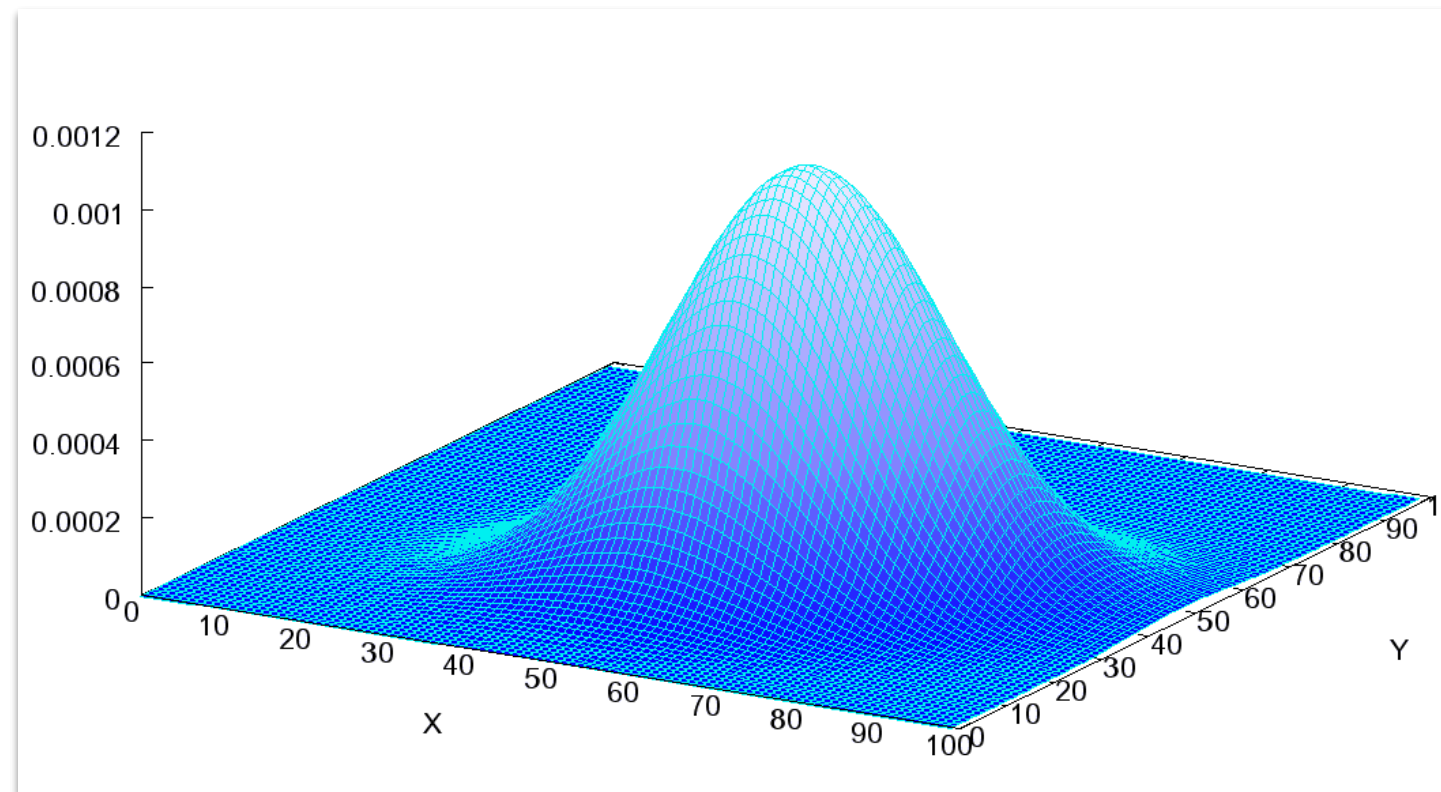
# Properties of Gaussians

# Normal Distribution



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

# Multivariate Normal



Vectors

$(X_1, \dots, X_D)$

$(\mu_1, \dots, \mu_D)$

Density:  $p(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^D |\boldsymbol{\Sigma}|}} \exp^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})}$

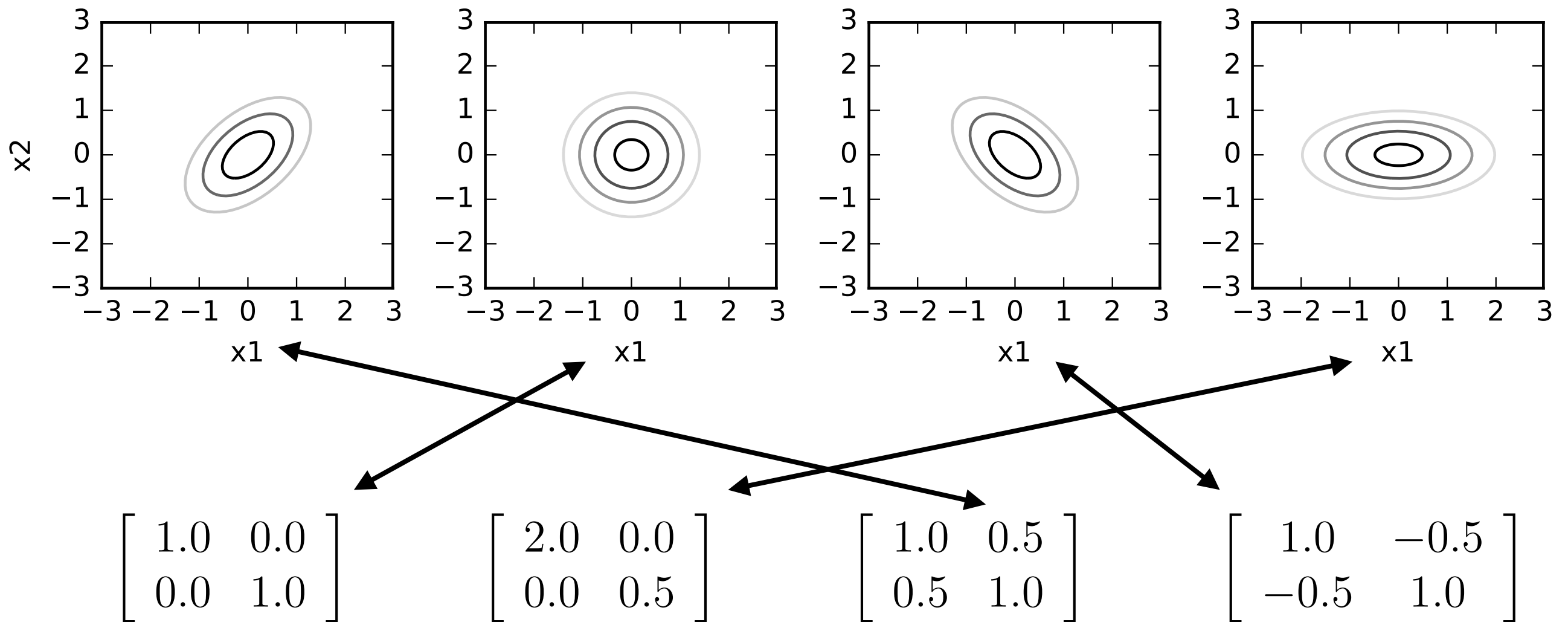
Covariance Matrix

Mean:  $\mathbb{E}[X_d] = \mu_d$

Covariance:  $\text{Cov}[X_d, X_e] = \Sigma_{de}$

# Covariance Matrices

Density:  $f(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^D |\boldsymbol{\Sigma}|}} \exp^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})}$



*Question:* Which covariance matrix  $\boldsymbol{\Sigma}$  corresponds to which plot?

# Marginals and Conditionals

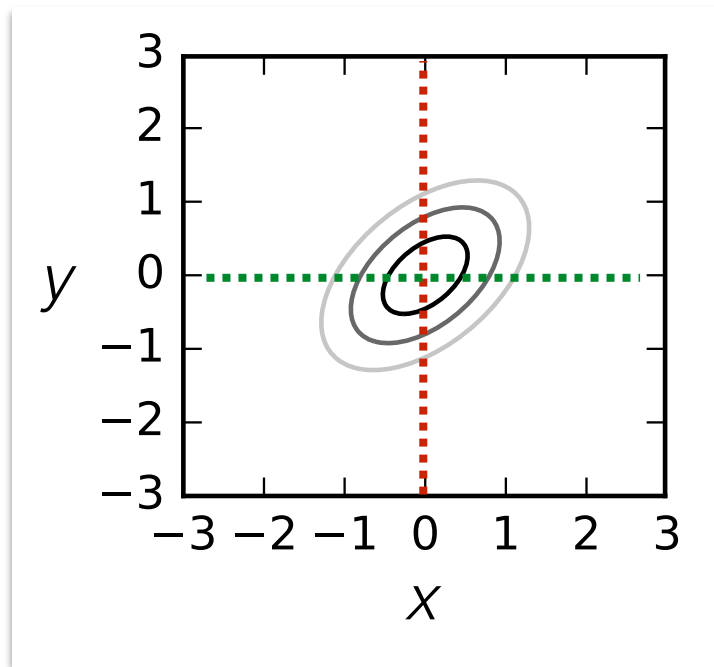
Suppose that  $\mathbf{x}$  and  $\mathbf{y}$  are jointly Gaussian:

$$\mathbf{z} = \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \sim \mathcal{N} \left( \begin{bmatrix} \mathbf{a} \\ \mathbf{b} \end{bmatrix}, \begin{bmatrix} \mathbf{A} & \mathbf{C} \\ \mathbf{C}^T & \mathbf{B} \end{bmatrix} \right)$$

*Question:* What are the marginal distributions  $p(\mathbf{x})$  and  $p(\mathbf{y})$ ?

$$\mathbf{x} \sim \mathcal{N}(\mathbf{a}, \mathbf{A})$$

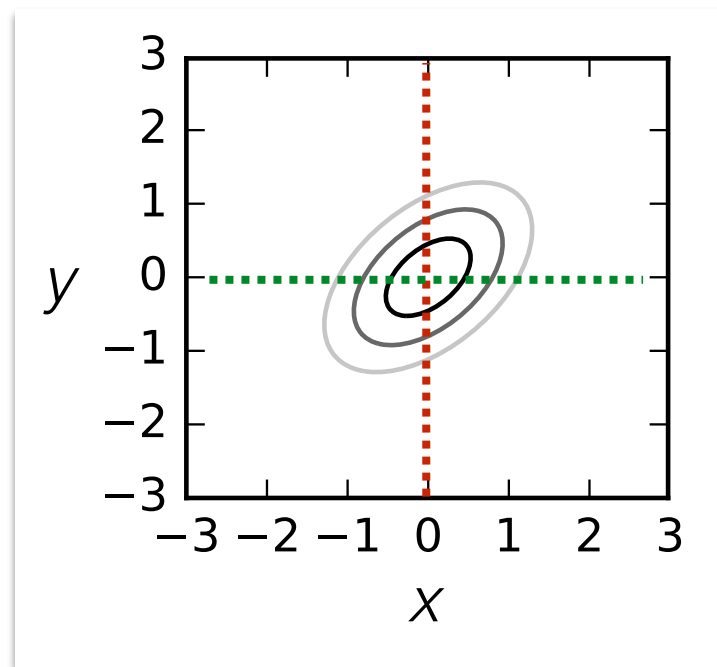
$$\mathbf{y} \sim \mathcal{N}(\mathbf{b}, \mathbf{B})$$



# Marginals and Conditionals

Suppose that  $\mathbf{x}$  and  $\mathbf{y}$  are jointly Gaussian:

$$\mathbf{z} = \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \sim \mathcal{N} \left( \begin{bmatrix} \mathbf{a} \\ \mathbf{b} \end{bmatrix}, \begin{bmatrix} \mathbf{A} & \mathbf{C} \\ \mathbf{C}^T & \mathbf{B} \end{bmatrix} \right)$$



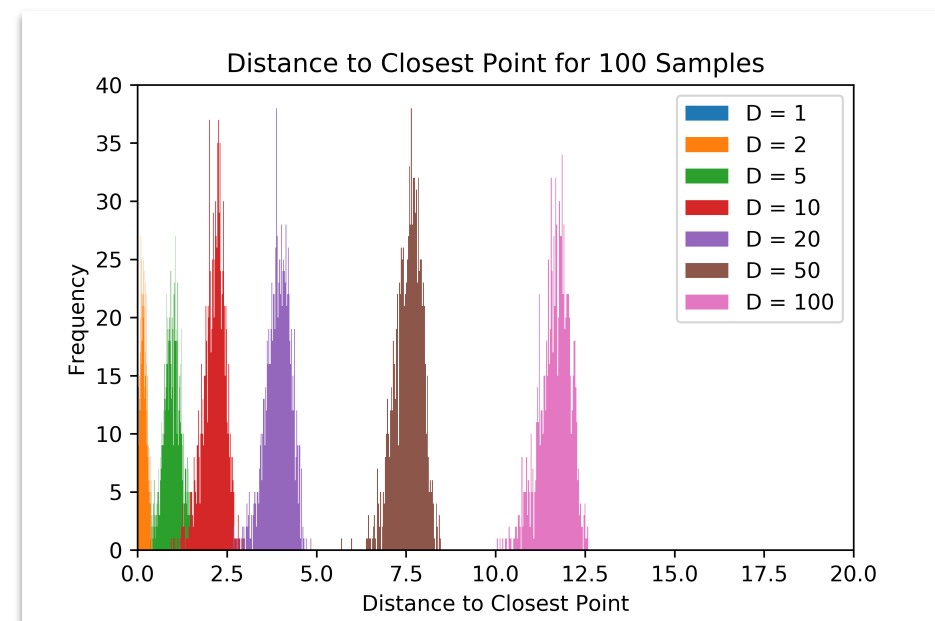
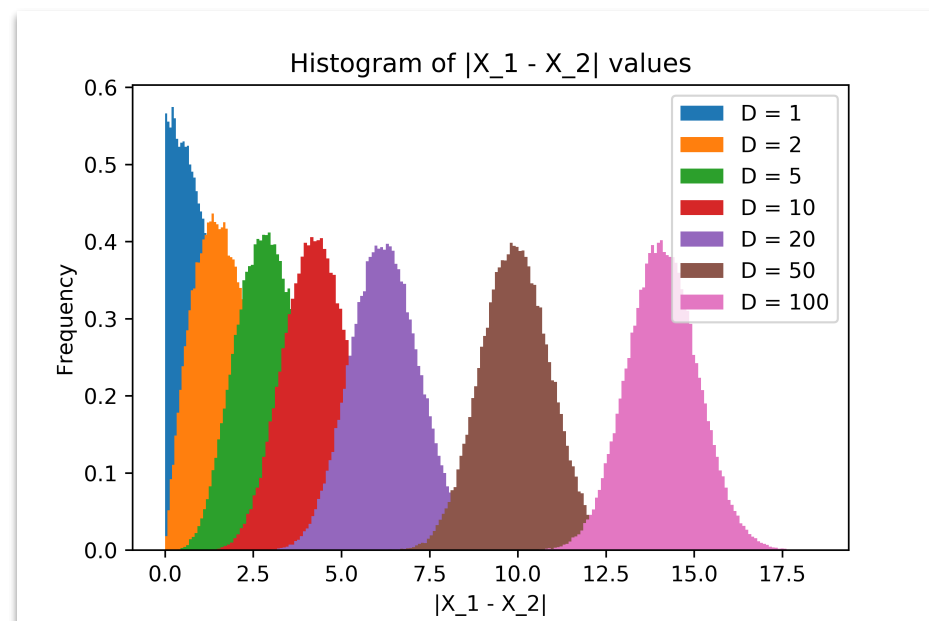
Once can derive the conditional distributions  $p(\mathbf{x} | \mathbf{y})$  and  $p(\mathbf{y} | \mathbf{x})$  in closed form as well; they are also Normals!

# Curse of Dimensionality

**Question:** Suppose that  $X_1$  and  $X_2$  are independent Gaussian variables with diagonal covariance

$$X_1 \sim \text{Normal}(0, \sigma^2 I_D) \quad X_2 \sim \text{Normal}(0, \sigma^2 I_D)$$

How does the distribution on the distance  $|X_1 - X_2|$  change as we increase the dimension  $D$ ?



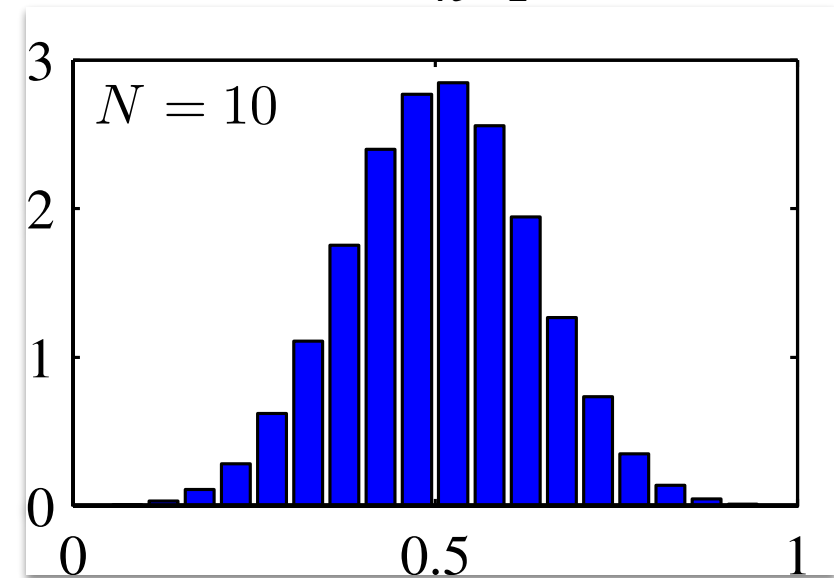
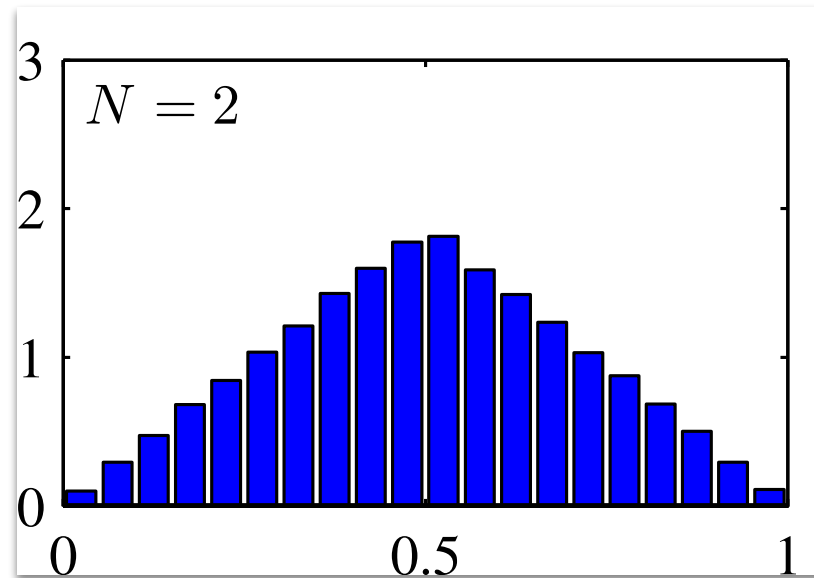
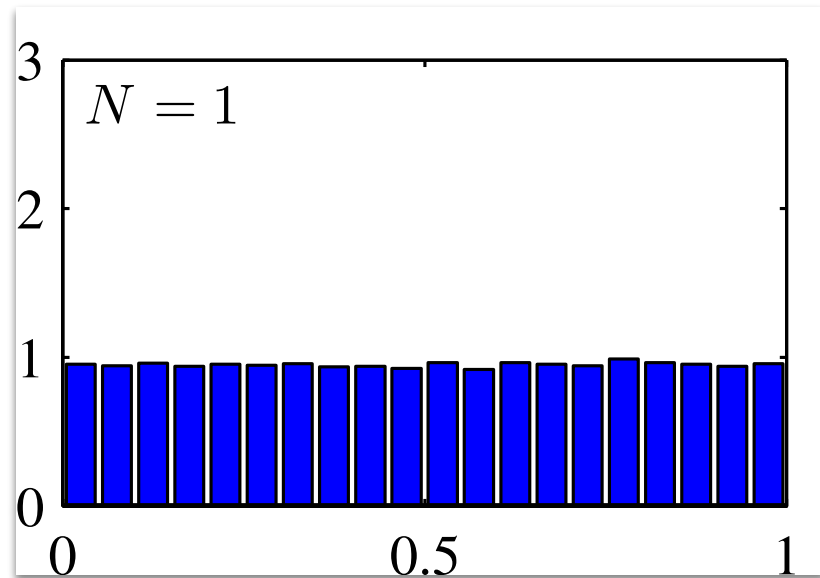


# Central Limit Theorem

$$\bar{X} = X_1$$

$$\bar{X} = \frac{1}{2} (X_1 + X_2)$$

$$\bar{X} = \frac{1}{N} \sum_{n=1}^N X_n$$



If  $X_1, \dots, X_N$  are

1. Independent identically distributed (i.i.d.)
2. Have finite variance  $0 < \sigma_X^2 < \infty$

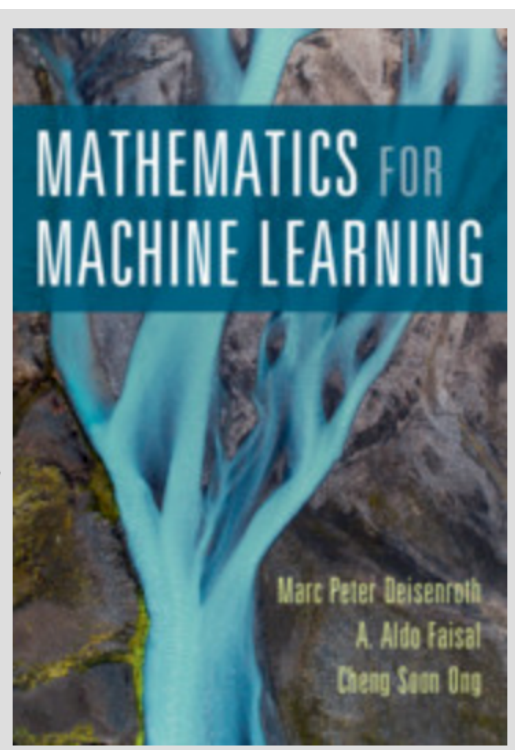
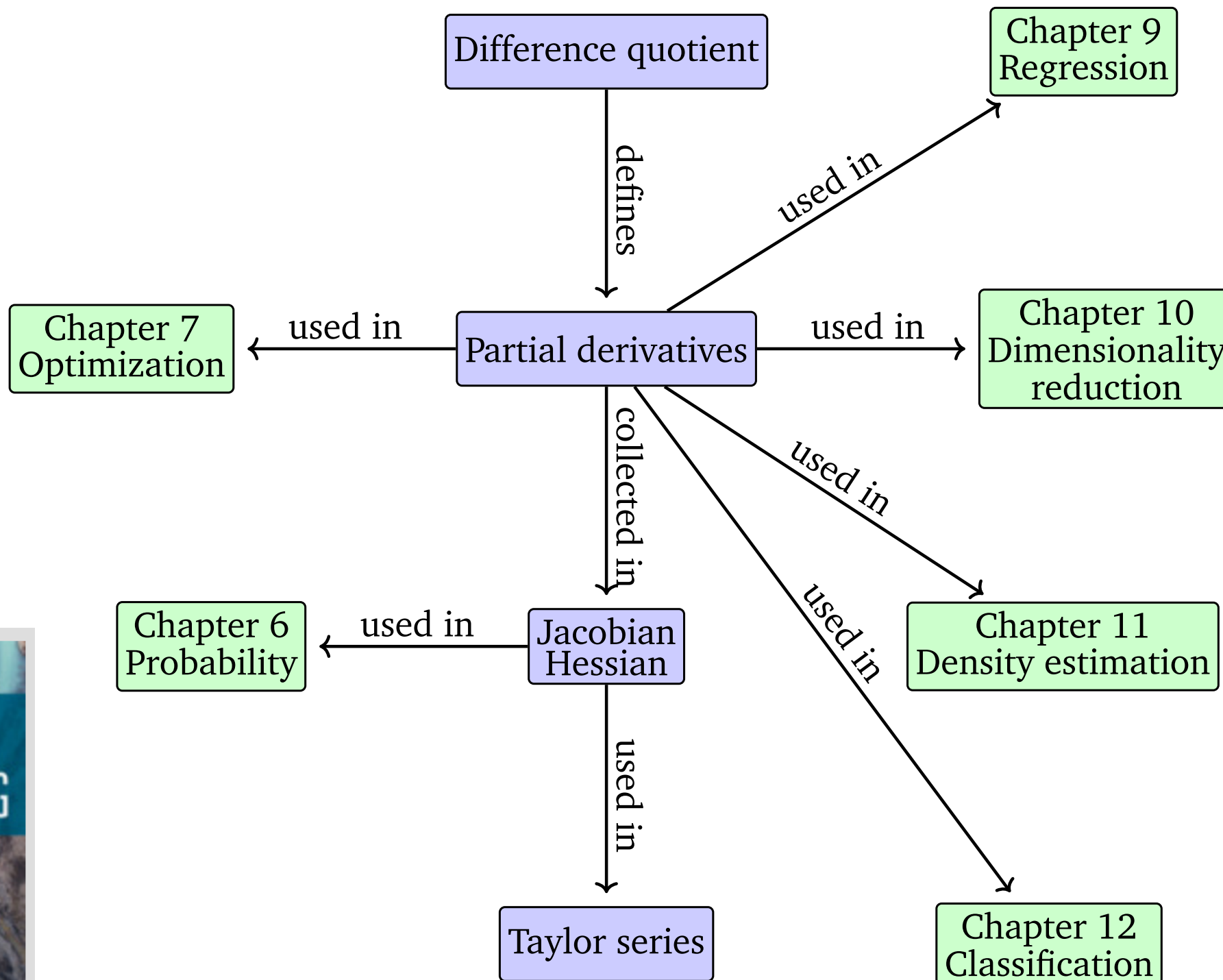
Then, as  $N$  approaches  $\infty$ , the mean is distributed as

$$p(\bar{x}) = \text{Normal}(\bar{x}; \mu_X, \sigma_X^2 / N) \quad \bar{X} = \frac{1}{N} \sum_{n=1}^N X_n$$

# Summary of Gaussians

In sum, the Normal (or Gaussian) distribution pops up everywhere and is easy to work with; familiarize yourself with it!

# Some Calc Review



# Univariate Functions

$$y = f(x), \quad x, y \in \mathbb{R}$$

*Difference Quotient*

$$\frac{\delta y}{\delta x} \stackrel{\text{def}}{=} \frac{f(x + \delta x) - f(x)}{\delta x}$$

# Univariate Functions

*Derivative (formally)*

$$\frac{df}{dx} := \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

# Univariate Functions

*Derivative (formally)*

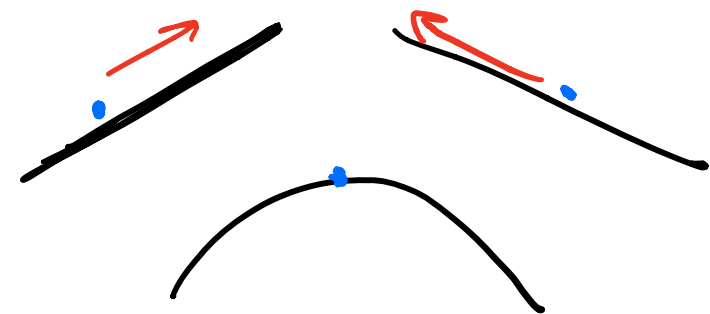
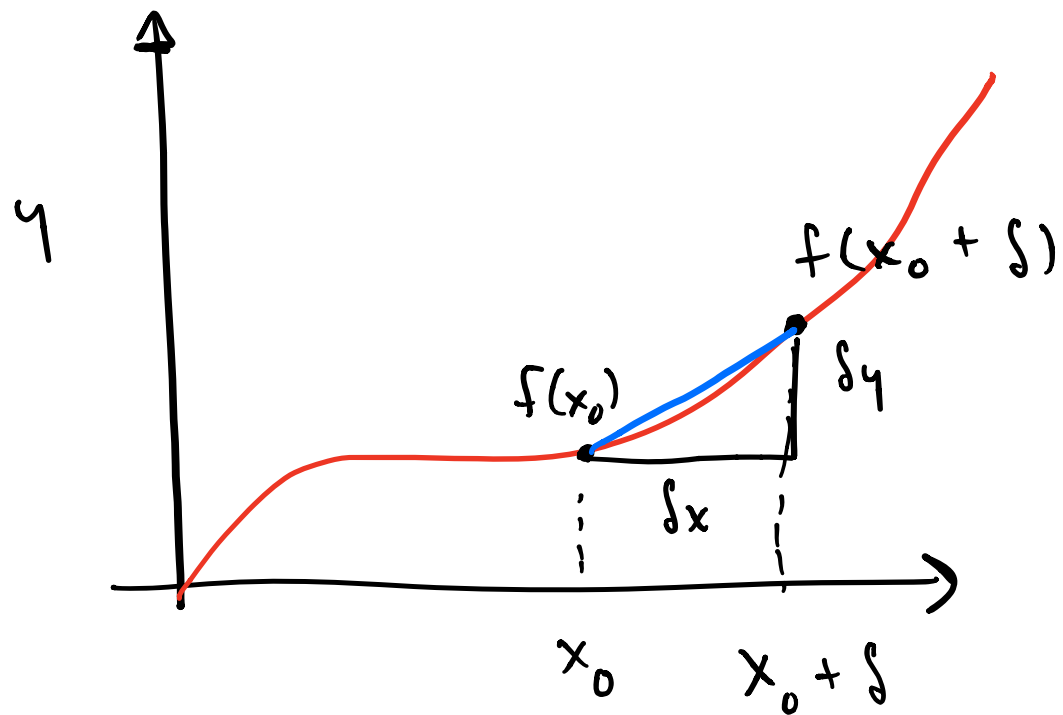
$$\frac{df}{dx} := \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

The derivative points in the direction of steepest ascent

Consider univariate case.

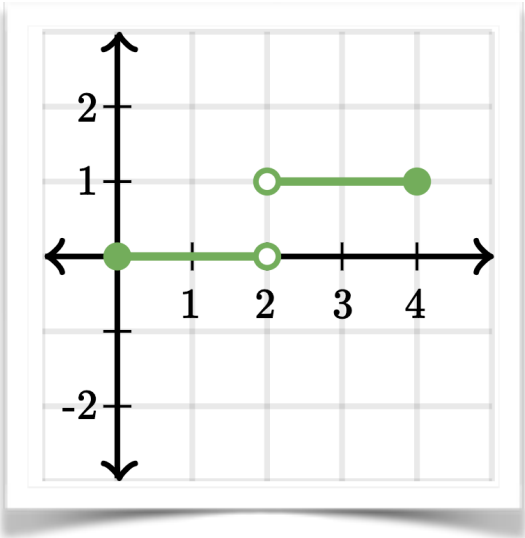
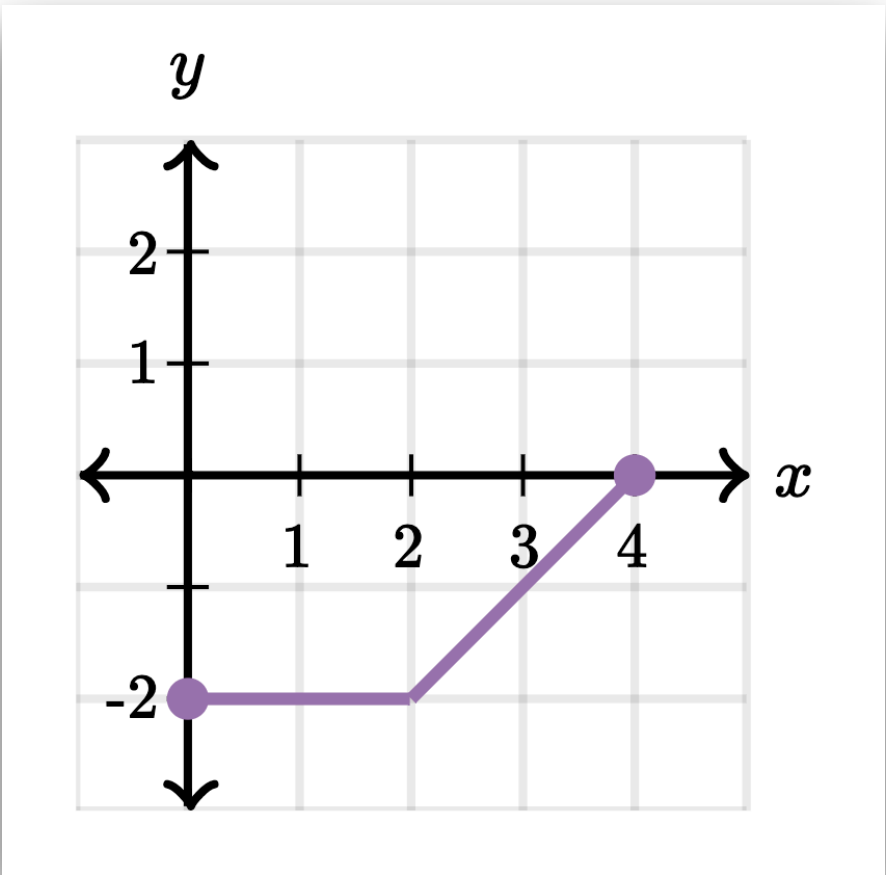
$$y = f(x)$$

$$\frac{\delta y}{\delta x} \stackrel{\text{def}}{=} \frac{f(x + \delta x) - f(x)}{\delta x}$$

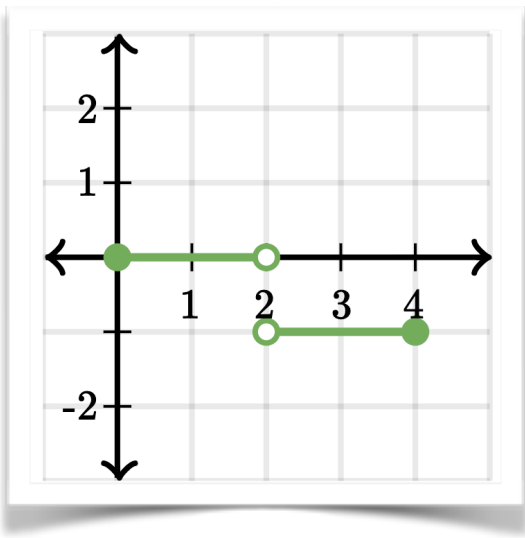




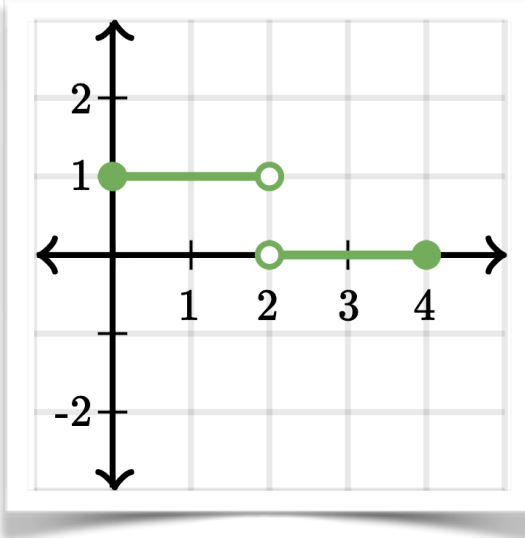
Spot the derivative



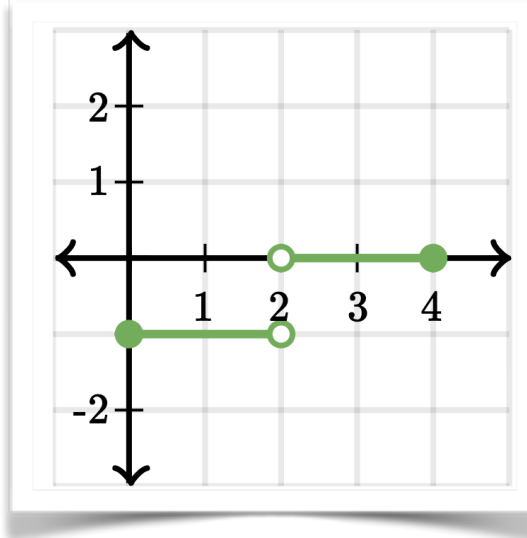
A



B

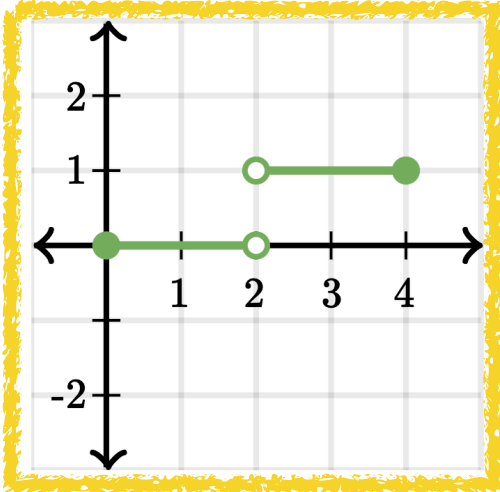
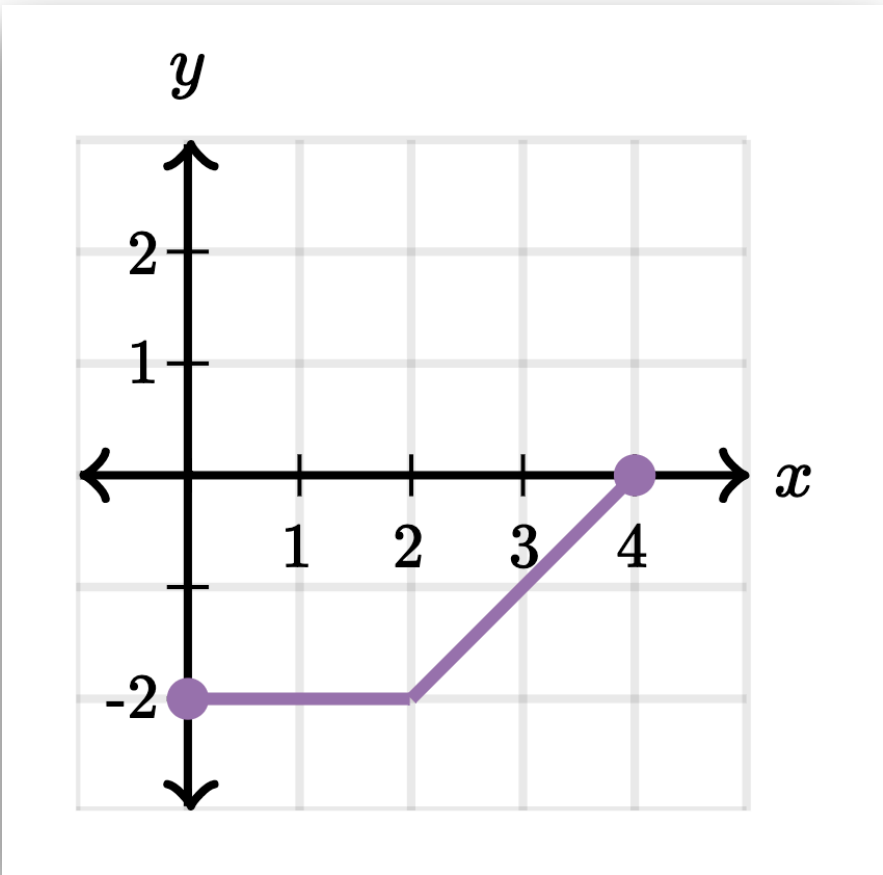


C

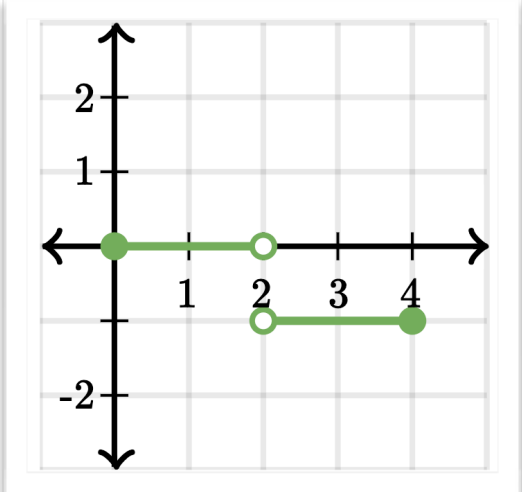


D

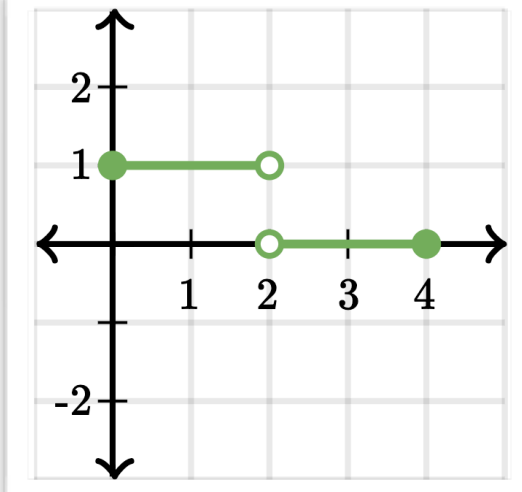
Spot the derivative



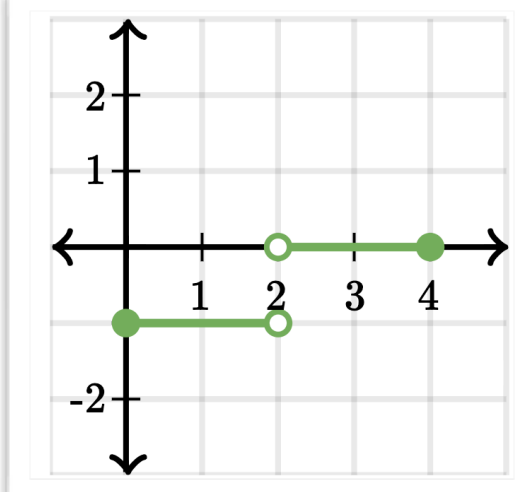
A



B



C



D

# Sum Rule

$$(f(x) + g(x))' = f'(x) + g'(x)$$

# Product Rule

$$(f(x)g(x))' = f'(x)g(x) + f(x)g'(x)$$

# Chain Rule

$$(g(f(x)))' = (g \circ f)'(x) = g'(f(x))f'(x)$$

# Gradients

Usually in ML we care about multivariate functions

$\boldsymbol{x} \in \mathbb{R}^n$  of  $n$  variables  $x_1, \dots, x_n$

$f : \mathbb{R}^n \rightarrow \mathbb{R}$

# Gradients

Partial derivatives are taken w.r.t. one dimension at a time:

$$\frac{\partial f}{\partial x_1} = \lim_{h \rightarrow 0} \frac{f(x_1 + h, x_2, \dots, x_n) - f(\mathbf{x})}{h}$$

⋮

$$\frac{\partial f}{\partial x_n} = \lim_{h \rightarrow 0} \frac{f(x_1, \dots, x_{n-1}, x_n + h) - f(\mathbf{x})}{h}$$

# Gradients

Group the gradients into a vector (the *gradient*)

$$\nabla_{\mathbf{x}} f = \text{grad } f = \frac{df}{d\mathbf{x}} = \left[ \frac{\partial f(\mathbf{x})}{\partial x_1} \quad \frac{\partial f(\mathbf{x})}{\partial x_2} \quad \dots \quad \frac{\partial f(\mathbf{x})}{\partial x_n} \right] \in \mathbb{R}^{1 \times n}$$



# Example

$$f(x_1, x_2) = x_1^2 x_2 + x_1 x_2^3$$

$$\frac{\partial f(x_1, x_2)}{\partial x_1} = 2x_1 x_2 + x_2^3$$

$$\frac{\partial f(x_1, x_2)}{\partial x_2} = x_1^2 + 3x_1 x_2^2$$

$$\frac{df}{d\mathbf{x}} = \left[ \frac{\partial f(x_1, x_2)}{\partial x_1} \quad \frac{\partial f(x_1, x_2)}{\partial x_2} \right] = [2x_1 x_2 + x_2^3 \quad x_1^2 + 3x_1 x_2^2] \in \mathbb{R}^{1 \times 2}$$

# Rules still hold!

Sum rule: 
$$\frac{\partial}{\partial \mathbf{x}} (f(\mathbf{x}) + g(\mathbf{x})) = \frac{\partial f}{\partial \mathbf{x}} + \frac{\partial g}{\partial \mathbf{x}}$$

Product rule: 
$$\frac{\partial}{\partial \mathbf{x}} (f(\mathbf{x})g(\mathbf{x})) = \frac{\partial f}{\partial \mathbf{x}} g(\mathbf{x}) + f(\mathbf{x}) \frac{\partial g}{\partial \mathbf{x}}$$

Chain rule: 
$$\frac{\partial}{\partial \mathbf{x}} (g \circ f)(\mathbf{x}) = \frac{\partial}{\partial \mathbf{x}} (g(f(\mathbf{x}))) = \frac{\partial g}{\partial f} \frac{\partial f}{\partial \mathbf{x}}$$

# Rules still hold!

Sum rule: 
$$\frac{\partial}{\partial \mathbf{x}} (f(\mathbf{x}) + g(\mathbf{x})) = \frac{\partial f}{\partial \mathbf{x}} + \frac{\partial g}{\partial \mathbf{x}}$$

Product rule: 
$$\frac{\partial}{\partial \mathbf{x}} (f(\mathbf{x})g(\mathbf{x})) = \frac{\partial f}{\partial \mathbf{x}} g(\mathbf{x}) + f(\mathbf{x}) \frac{\partial g}{\partial \mathbf{x}}$$

Chain rule: 
$$\frac{\partial}{\partial \mathbf{x}} (g \circ f)(\mathbf{x}) = \frac{\partial}{\partial \mathbf{x}} (g(f(\mathbf{x}))) = \frac{\partial g}{\partial f} \frac{\partial f}{\partial \mathbf{x}}$$

... but be mindful of dims!

## For review: Problem 5.7 in MML

5.7 Compute the derivatives  $df/d\mathbf{x}$  of the following functions by using the chain rule. Provide the dimensions of every single partial derivative. Describe your steps in detail.

a.

$$f(z) = \log(1 + z), \quad z = \mathbf{x}^\top \mathbf{x}, \quad \mathbf{x} \in \mathbb{R}^D$$

b.

$$f(\mathbf{z}) = \sin(\mathbf{z}), \quad \mathbf{z} = \mathbf{A}\mathbf{x} + \mathbf{b}, \quad \mathbf{A} \in \mathbb{R}^{E \times D}, \mathbf{x} \in \mathbb{R}^D, \mathbf{b} \in \mathbb{R}^E$$

where  $\sin(\cdot)$  is applied to every element of  $\mathbf{z}$ .

$$(a) \quad f(z) = \log(\overbrace{1+z}^u) \quad z = x^T x \quad x \in \mathbb{R}^D$$

$$\frac{df}{dx} = \frac{d}{du} \log u \frac{d}{dx} (1 + x^T x)$$

$$= \frac{1}{(1+x^T x)} 2x = \frac{2x}{(1+x^T x)}$$

$$(b) \quad f(z) = \sin(z) \quad z = Ax + b \quad A \in \mathbb{R}^{E \times D} \quad x \in \mathbb{R}^D$$

$$b \in \mathbb{R}^E$$

$$\frac{\partial f}{\partial x} = \frac{d \sin(m)}{dm} \frac{dm}{dx}$$

$$= \cos(m) \frac{\partial}{\partial x} Ax + b$$

$$= \underbrace{\cos(Ax + b)}_{E \times 1} \cdot \underbrace{A}_{E \times D}$$

intermezzo: the joys of auto-diff...

... or, first steps in pytorch