

# Important Announcements and Notes (Sep 19)

- Outcome vs. event
- Typo in general formula for marginal distribution

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

$$\text{outcome} = 3 \quad X = 3$$

$$\text{event} = \{4, 5, 6\}$$

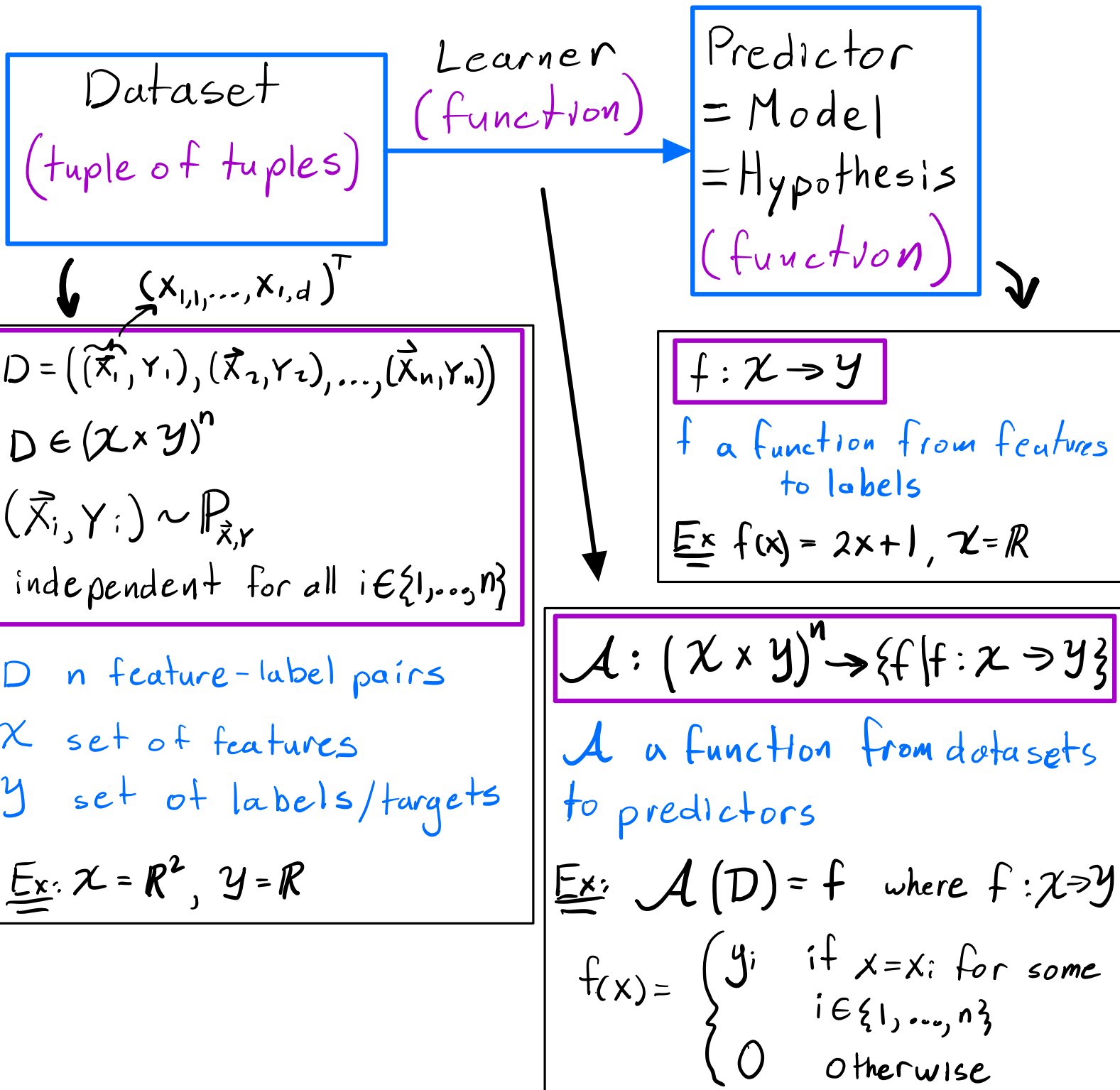
$$P(\{4, 5, 6\}) = P(X \in \{4, 5, 6\})$$

# Important Announcements and Notes (Sep 17)

- The probability distribution  $P$  is the thing you really care about (inputs are events)
  - You should think of the pmf and pdf (inputs are outcomes) as functions that help you calculate  $P$
- Countable set: A set with cardinality that is finite or countably infinite
- Uncountable set: A set with cardinality that is uncountably infinite
- No more discrete or continuous outcome space
- Review discrete and continuous r.v.
- A r.v. does not need to take values in  $\mathbb{R}$

# Motivation

Supervised Learning: Learning from a randomly sampled batch of labeled data



# Probability

Notes: Humans have a bad intuition when it comes to randomness

- Thinking Fast and Slow  
by: Daniel Kahneman

- Random variables
- Calculating probabilities using pmf and pdf
- Multivariate random variables
  - Conditional and marginal probabilities
- Representing random features, labels, and datasets
- Functions of random variables
- Expectation and variance

Warning: If some things seem informal, it is likely because we would need tools from Measure Theory, which we will not cover in this course.

Experiment: A process that generates an uncertain outcome

Ex: flipping a coin, rolling a dice

Outcome Space/Set: The set of all outcomes from the experiment

Ex:  $y = \{0, 1\}$  <sup>Heads</sup> <sub>Tails</sub> flipping a coin

$X = \{1, 2, 3, 4, 5, 6\}$  rolling a dice

$[0, 100]$

amount of a chemical in a wine

$\mathbb{R}$

measurement error

Discrete Outcome Space: finite or countably infinite

Ex:  $Y, X, \mathbb{N}$

Continuous Outcome space: uncountably infinite

Ex:  $[0, 100], \mathbb{R}$

Event: A subset of the outcome space (imprecise)

Ex: Outcome space:  $Y = \{0, 1\}$

Events:  $\{0\}, \{1\}, \{0, 1\} = Y, \emptyset$

Ex: Outcome space:  $X = \{1, 2, 3, 4, 5, 6\}$

Events:  $\emptyset, X, \{1\}, \{1, 3, 5\}, \{2, 4, 6\}, \{1, 2, 3\}, \dots$

Ex: Outcome space:  $[0, 100]$

Events:  $\emptyset, [0, 100], [0, 4], [1, 2] \cup [7, 30], \dots$

Probability Distribution: A function  $P$  defining the likelihood of each event (and satisfying certain properties)

$P: \underbrace{\text{event space/set}} \rightarrow [0, 1]$  Think of this as a set containing all the events

A complicated set ( $\sigma$ -algebra) that we will not define

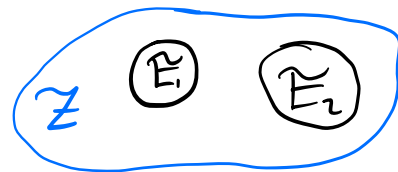
Properties: (imprecise)

Outcome space:  $\mathcal{Z}$

1.  $P(\mathcal{Z}) = 1$

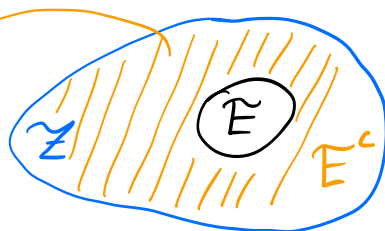
2. If  $E_1 \subset \mathcal{Z}, E_2 \subset \mathcal{Z}$  and  $E_1 \cap E_2 = \emptyset$ , then

$$P(E_1 \cup E_2) = P(E_1) + P(E_2)$$



Ex: (of property 2.)

Events:  $E, E^c$



$$E \cap E^c = \emptyset, E \cup E^c = \mathcal{Z}$$

$$P(E \cup E^c) = P(E) + P(E^c)$$

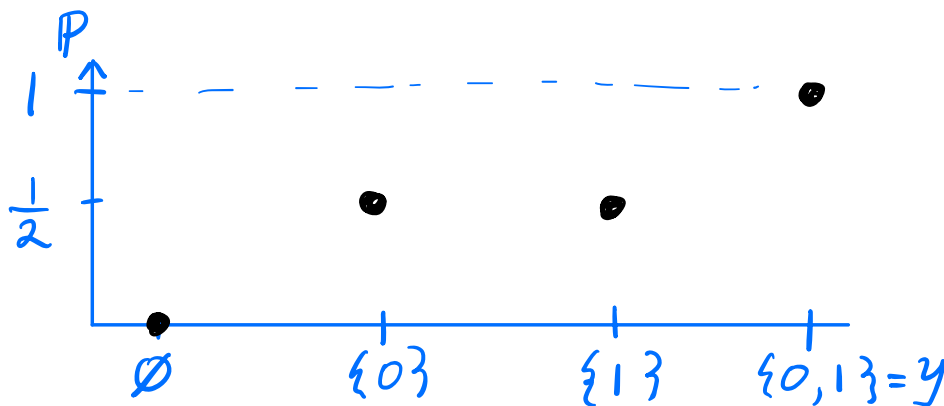
$$= P(\mathcal{Z}) = 1$$

rearranging:

$$P(E) = 1 - P(E^c)$$

Ex: Outcome space:  $\mathcal{Y} = \{0, 1\}$

$$P(\emptyset) = 0, P(\mathcal{Y}) = 1, P(\{0\}) = \frac{1}{2}, P(\{1\}) = \frac{1}{2}$$

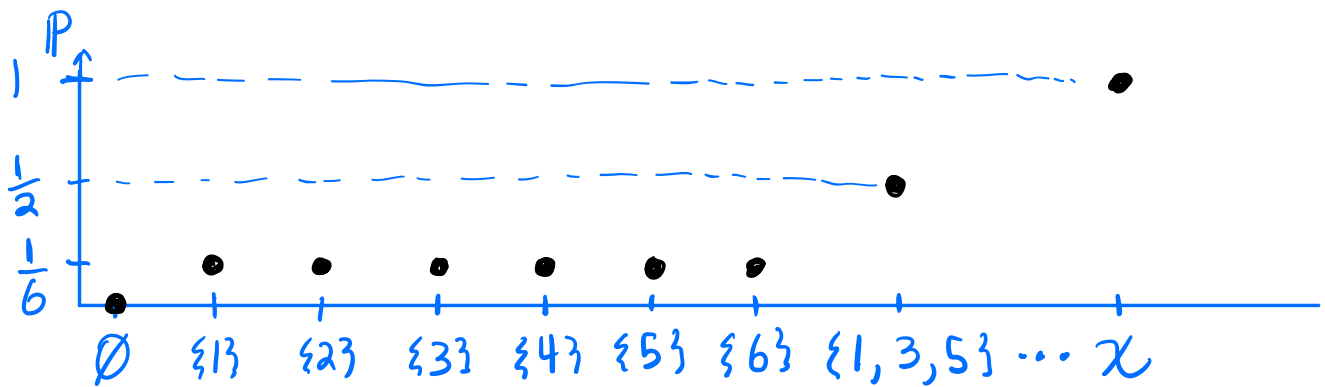


Ex: Outcome space:  $\mathcal{X} = \{1, 2, 3, 4, 5, 6\}$

$$P(\{1\}) = P(\{2\}) = \dots = P(\{6\}) = \frac{1}{6}$$

$$P(\{1, 3, 5\}) = P(\{1\}) + P(\{3\}) + P(\{5\}) = \frac{1}{2}$$

$$P(\mathcal{X}) = 1$$



## Random Variables

Random Variable (r.v.): (imprecise) A variable that takes a value based on the outcome of an experiment, and is associated with a probability distribution. ~~The value is always in  $\mathbb{R}$~~

Ex:  $X \in \mathcal{X} = \{1, 2, 3, 4, 5, 6\}$  with IP from prev. example  
 $Y \in \mathcal{Y} \in \{0, 1\}$  with IP from prev. example  
 ~~$Z \in \{H, T\}$  is not a r.v. since  $H \notin \mathbb{R}, T \notin \mathbb{R}$~~

A random variable is actually a function (satisfying certain properties) from one outcome space to another outcome space. Ex  $X(T) = 0, X(H) = 1$ .

\*It will not be necessary to know this for this course\*

# Probability Distributions with r.v.

Ex: Outcome space:  $\mathcal{X}$  r.v.:  $X \in \mathcal{X}$

$$P(\{1, 3, 5\}) \stackrel{\text{def}}{=} P(X \in \{1, 3, 5\})$$

$$P(\{4, 5, 6\}) = P(X \in \{4, 5, 6\}) = P(X \geq 4)$$

$$P(\{4\}) = P(X \in \{4\}) = P(X=4)$$

Notation:  $Z \sim P$  "Z is sampled according to distribution P"

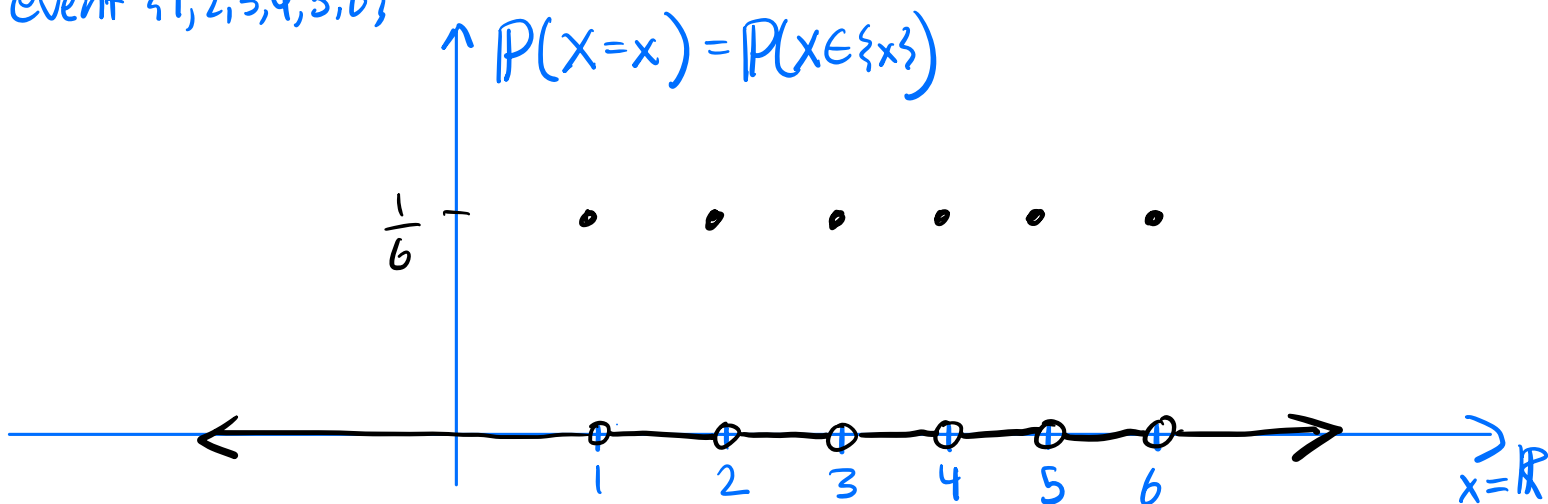
Discrete r.v.: A r.v. that takes values from:

- A countable outcome space, or
- an uncountable outcome space, but there is a countable event that has probability 1

Ex:  $Y \in \mathcal{Y} = \{0, 1\}$ ,  $X \in \mathcal{X} = \{1, 2, 3, 4, 5, 6\}$ ,  $Z \in \mathbb{N}$

Ex:  $X \in \mathbb{R}$  where  $P(X=1) = \dots = P(X=6) = \frac{1}{6}$

Probability 1  $\xrightarrow{\text{so}}$   $P(X \in \{1, 2, 3, 4, 5, 6\}) = 1$   
for countable event  $\{1, 2, 3, 4, 5, 6\}$  and  $P(\mathbb{R} \setminus \{1, 2, 3, 4, 5, 6\}) = 0$





Note: You can always take a r.v. defined on a countable outcome space and define it on a larger uncountable outcome space by setting the probability of the event containing all the new outcomes to zero

Continuous r.v.: A r.v. that takes values from:

- an uncountable outcome space and the probability of any single outcome is zero

Ex:  $Z \in [0, 900]$  and  $P(Z=z) = P(Z \in \{z\}) = 0$  for all  $z \in [0, 900]$   
but  $P(Z \in [0, 900]) = 1$

Ex:  $Z \in \mathbb{R}$  and  $P(Z=z) = P(Z \in \{z\}) = 0$  for all  $z \in \mathbb{R}$   
but  $P(Z \in \mathbb{R}) = 1$

## Calculating Probabilities

Motivation: It is hard to define the values of a probability distribution  $P$  for all the events

Probability Mass Function (pmf): A function  $p: \mathcal{Z} \rightarrow [0, 1]$

where  $\mathcal{Z}$  is a discrete outcome space and  $\sum_{z \in \mathcal{Z}} p(z) = 1$ .

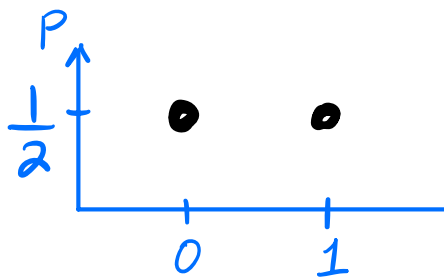
The probability of an event  $E \subset \mathcal{Z}$  is:

$$P(Z \in E) \stackrel{\text{def}}{=} \sum_{z \in E} p(z)$$

where  $Z \in \mathcal{Z}$

Ex: Outcome space:  $\mathcal{Y}$

$$p(0) = \frac{1}{2}, p(1) = \frac{1}{2}$$



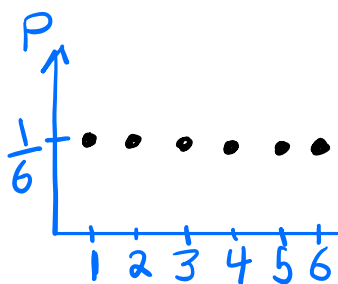
$$P(Y \in \{0, 1\}) = \sum_{Y \in \{0, 1\}} p(y) = p(0) + p(1) = 1$$

$$P(Y=0) = P(Y=1) = P(Y \in \{1\}) = \sum_{Y \in \{1\}} p(y) = p(1) = \frac{1}{2}$$

$$P(Y \in \emptyset) = \sum_{Y \in \emptyset} p(y) = 0$$

Ex: Outcome space:  $\mathcal{X}$

$$p(1) = p(2) = \dots = p(6) = \frac{1}{6}$$



$$P(X \in \{1, 3, 5\}) = \sum_{X \in \{1, 3, 5\}} p(x) = p(1) + p(3) + p(5) = \frac{1}{2}$$

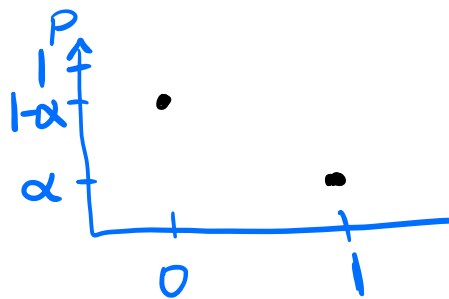
Discrete Probability Distributions with special names:

Bernoulli distribution (parameter:  $\alpha \in [0, 1]$ ):

Outcome space:  $\{0, 1\}$

pmf:  $p(1) = \alpha, p(0) = 1 - \alpha$

Distribution  $P = \text{Bernoulli}(\alpha)$



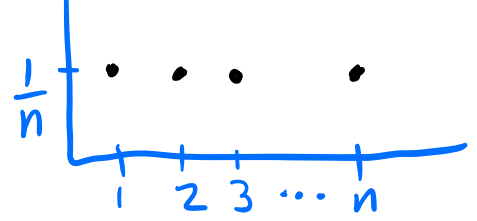
$P(Z=1) = P(Z \in \{1\}) = p(1) = \alpha$   $Z \in \{0, 1\}$  is a "Bernoulli r.v."

Discrete Uniform Distribution (parameter:  $n$ ):

Outcome space:  $\{1, 2, \dots, n\}$

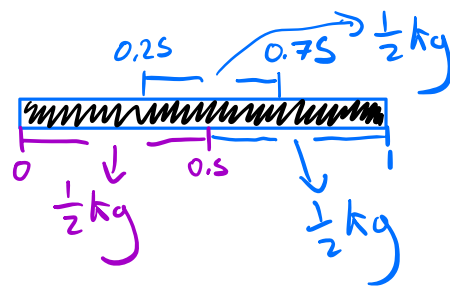
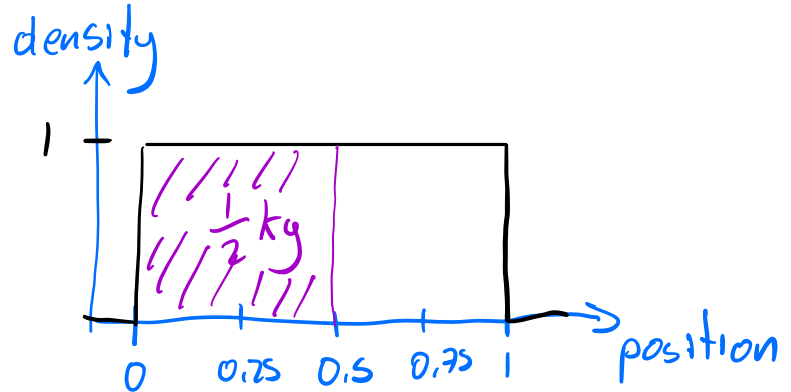
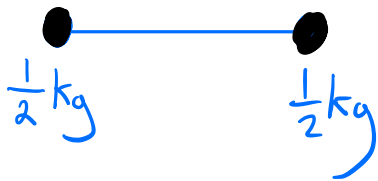
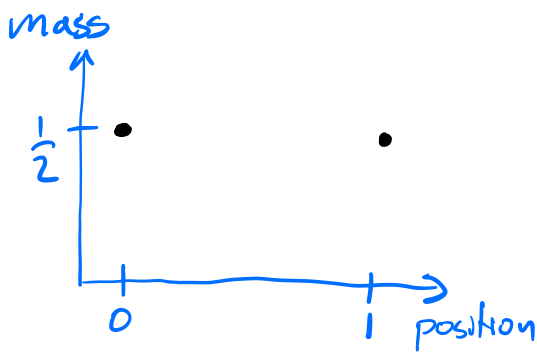
$P \uparrow$

pmf:  $p(1) = p(2) = \dots = p(n) = \frac{1}{n}$



Distribution  $P = \text{Uniform}(n)$

## Intuition with a rod in physics



Probability Density Function (pdf): a function  $p: \mathcal{Z} \rightarrow [0, \infty]$

where  $\mathcal{Z}$  is an **uncountable** outcome space and  $\int_{\mathcal{Z}} p(z) dz = 1$

The probability of an event  $E \subset \mathcal{Z}$  is:

$$P(Z \in E) \stackrel{\text{def}}{=} \int_E p(z) dz$$

$$P(Z = z) = P(Z \in \{z\}) = 0 \quad \text{where } z \in \mathcal{Z}$$

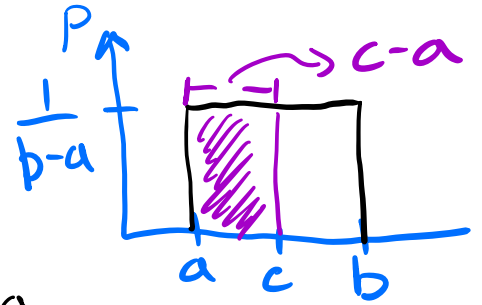
# Continuous Probability Distributions with special names:

## Continuous Uniform Distribution (parameters: $a \in \mathbb{R}, b \in \mathbb{R}$ ):

Outcome space:  $[a, b]$

$$\text{pdf: } p(z) = \frac{1}{b-a}$$

Distribution  $P = \text{Uniform}(a, b)$



$$P(a \leq Z \leq c) = P(Z \in [a, c]) = \int_a^c p(z) dz = \frac{c-a}{b-a}$$

where  $a \leq c \leq b$

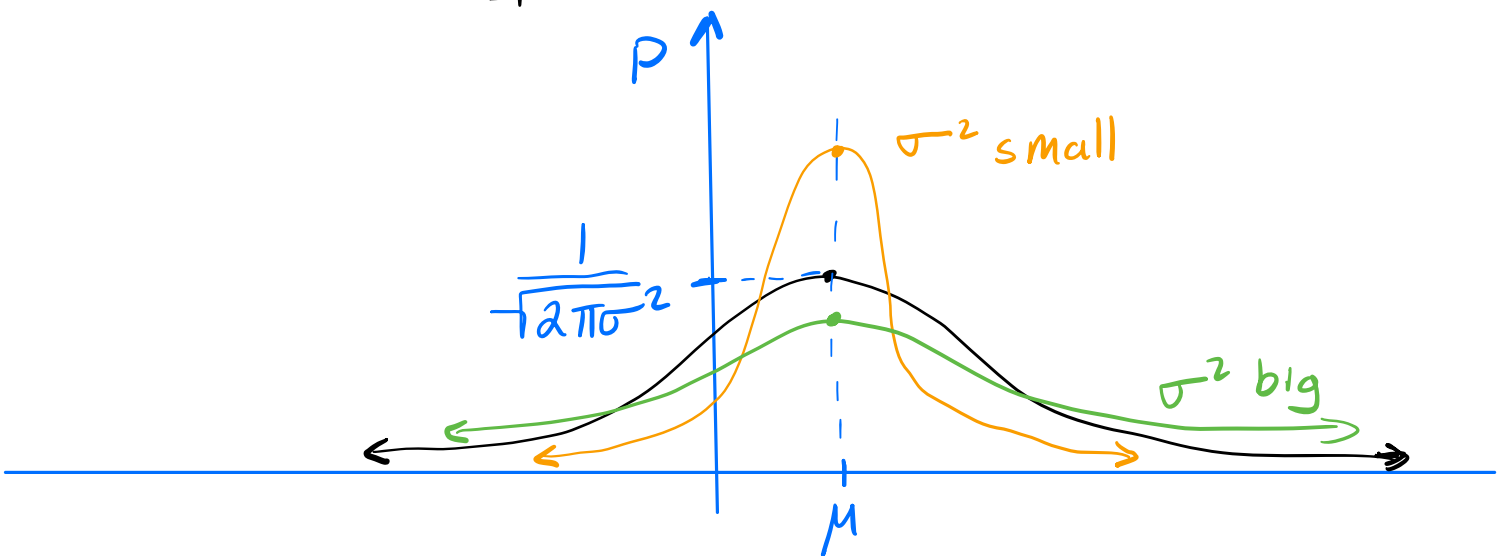
## Gaussian/Normal Distribution (parameters: $\mu \in \mathbb{R}, \sigma^2 > 0$ ):

Outcome space:  $\mathbb{R}$

$$\text{pdf: } p(z) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(z-\mu)^2\right)$$

Distribution  $P = \mathcal{N}(\mu, \sigma^2) = \text{Gaussian}(\mu, \sigma^2)$

$$\underline{\underline{\text{Ex}}} \quad P(-1 \leq Z \leq 1) = \int_{-1}^1 p(z) dz$$

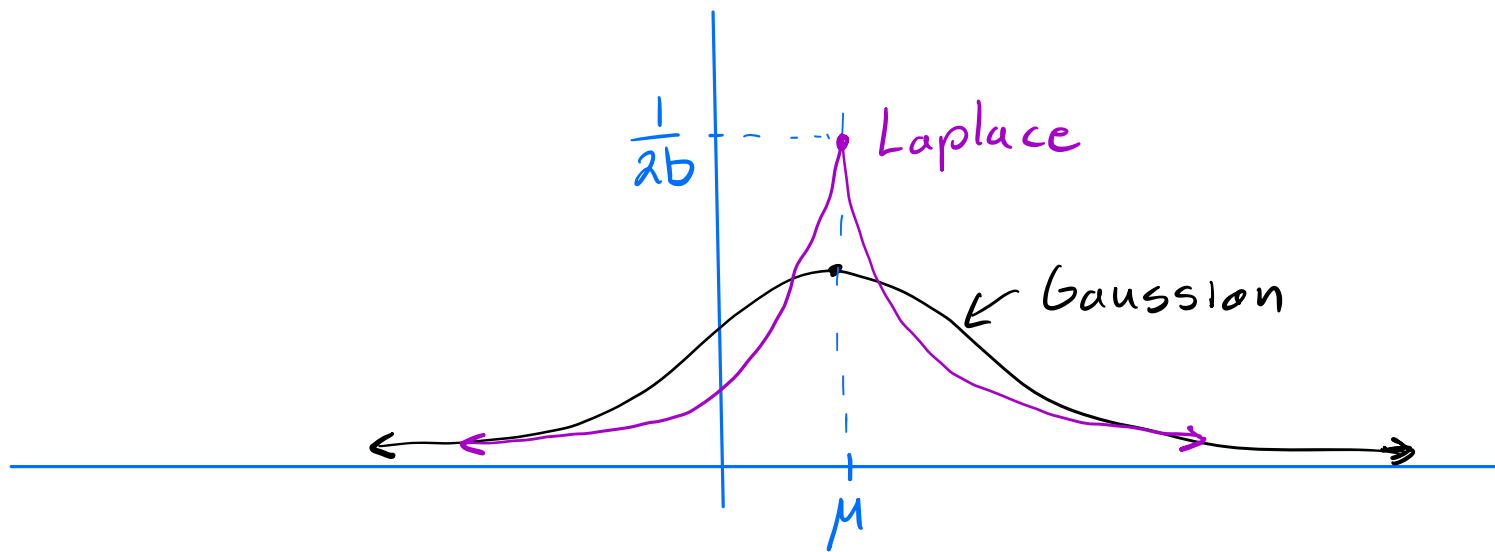


## Laplace Distribution (parameters: $\mu \in \mathbb{R}, b > 0$ ):

Outcome space:  $\mathbb{R}$

$$\text{pdf: } p(x) = \frac{1}{2b} \exp\left(-\frac{1}{b}|x - \mu|\right)$$

$$\text{Distribution } \mathbb{P} = \text{Laplace}(\mu, b)$$



## Multivariate Random Variables

Motivation: To be able to talk about the probability of different types of events at the same time

Ex: The probability of getting heads and rolling a 3

The probability of a wine containing 2.5mg of one chemical and 4mg of another chemical

The probability of a house having 4 rooms and 2 washrooms and being less than 10min from a university

The probability of being young and having arthritis

Multi variate Random Variable: A tuple of more than one random variable

Ex: (Flipping 2 coins) Heads

Outcome space:  $\mathcal{X} = \{0,1\} \times \{0,1\} = \{(0,0), (0,1), (1,0), (1,1)\}$

r.v.:  $X = (X_1, X_2) \in \mathcal{X}$

(Collecting the info of one house (ex: # of rooms, age))

Outcome space:  $\mathcal{X} = \mathbb{N} \times [0, \infty)$  age

r.v.:  $X = (X_1, X_2) \in \mathcal{X}$  # of rooms

(Collecting the info of one house and its price)

Outcome space:  $\mathcal{Z} = (\mathbb{N} \times [0, \infty)) \times [0, \infty)$  age

r.v.:  $Z = (X, Y)$   
 $= ((X_1, X_2), Y)$  # of rooms Price

Calculating Joint Probabilities

Ex: (If you have arthritis and if you are young or old)

Outcome space:  $\mathcal{Z} = \{0,1\} \times \{0,1\}$  Old Arthritis  
Young No arthritis

r.v.:  $Z = (X, Y) \in \mathcal{Z}$

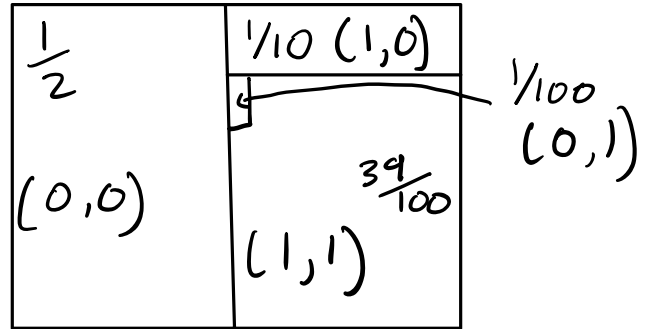
$= \{(0,0), (0,1), (1,0), (1,1)\}$

pmf:  $p: \mathcal{Z} \rightarrow [0, 1]$

$\mathcal{Z} = \{(0,0), (0,1), (1,0), (1,1)\}$

Not based on real data  $\left\{ \begin{array}{l} p((0,0)) = p(0,0) = \frac{1}{2}, \quad p(0,1) = \frac{1}{100} \\ p(1,0) = \frac{1}{10}, \quad p(1,1) = \frac{39}{100} \end{array} \right.$

|   |   | Y              |                  |
|---|---|----------------|------------------|
|   |   | 0              | 1                |
| X | 0 | $\frac{1}{2}$  | $\frac{1}{100}$  |
|   | 1 | $\frac{1}{10}$ | $\frac{39}{100}$ |



$$\sum_{z \in \mathcal{Z}} p(z) = \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x,y) = \frac{1}{2} + \frac{1}{100} + \frac{1}{10} + \frac{39}{100} = 1$$

What is the probability of being young (i.e.  $X=0$ )?

$$\mathcal{E} = \{(0,0), (0,1)\} = \{0\} \times \{0,1\} \subset \mathcal{Z}$$

$$\begin{aligned} P(X=0, Y \in \{0,1\}) &= P(Z \in \mathcal{E}) = \sum_{z \in \mathcal{E}} p(z) \\ &= \sum_{x \in \{0\}} \sum_{y \in \{0,1\}} p(x,y) \\ &= p(0,0) + p(0,1) \\ &= \frac{1}{2} + \frac{1}{100} = \frac{51}{100} \end{aligned}$$

Marginal Distribution: The distribution over a subset of random variables

$\mathbb{E}_x$ : Continuing with the arthritis example

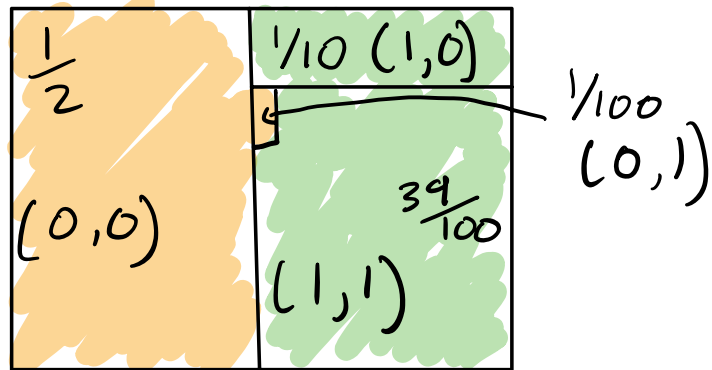
Marginal Distribution:  $P_x (X \in \tilde{E}_x)$  where  $\tilde{E}_x \in \mathcal{X}$

Marginal pmf:  $p_x: \mathcal{X} \rightarrow [0, 1]$ ,  $p_x(x) = \sum_{y \in \mathcal{Y}} p(x, y)$

$$P_x(X=0) \stackrel{P(X \in \{0\})}{=} p_x(0) = \sum_{x \in \{0\}} \sum_{y \in \mathcal{Y}} p(x, y) = \frac{51}{100}$$

$$P_x(X=1)$$

$$= \frac{1}{10} + \frac{39}{100} = \frac{49}{100}$$



Discrete r.v.  $X_1, \dots, X_d$

$X = (X_1, \dots, X_d) \in \mathcal{X}_1 \times \dots \times \mathcal{X}_d = \mathcal{X}$ ,  $p: \mathcal{X} \rightarrow [0, 1]$

$p_{x_i}: \mathcal{X}_i \rightarrow [0, 1]$ ,  $i \in \{1, \dots, d\}$

Marginal pmf:

$$p_{x_i}(x_i) \stackrel{\text{def}}{=} \sum_{x_1 \in \mathcal{X}_1} \dots \sum_{x_{i-1} \in \mathcal{X}_{i-1}} \sum_{x_{i+1} \in \mathcal{X}_{i+1}} \dots \sum_{x_d \in \mathcal{X}_d} p(x_1, \dots, x_d)$$

$$P_{x_i}(X_i \in \tilde{E}_i) = \sum_{x_i \in \tilde{E}_i} p_{x_i}(x_i)$$

where  $\tilde{E}_i \subset \mathcal{X}_i$

Continuous r.v.  $X_1, \dots, X_d$

$X = (X_1, \dots, X_d) \in \mathcal{X}_1 \times \dots \times \mathcal{X}_d = \mathcal{X}$ ,  $p: \mathcal{X} \rightarrow [0, \infty)$



$$P_{X_i}: \mathcal{X}_i \Rightarrow [0, \infty), \quad i \in \{1, \dots, d\}$$

Marginal pdf:

$$P_{X_i}(X_i) = \int \dots \int_{\mathcal{X}_1} \int \dots \int_{\mathcal{X}_{i-1}} \int \dots \int_{\mathcal{X}_{i+1}} \int \dots \int_{\mathcal{X}_d} P(X_1, \dots, X_d) dx_1 \dots dx_{i-1} dx_{i+1} \dots dx_d$$

Distribution:

$$P_{X_i}(X_i \in \tilde{E}_i) = \int_{\tilde{E}_i} P_{X_i}(X_i) dx_i \quad \text{where } \tilde{E}_i \subset \mathcal{X}_i$$

Conditional Distribution: Probability of a r.v. given info about another r.v.

Ex: Probability that I have arthritis given I am young

Let r.v. =  $Y \in \mathcal{Y}, X \in \mathcal{X}$

Discrete  $Y$  for any  $x \in \mathcal{X}$  that  $P_X(x) \neq 0$

$$P_{Y|X=x}: \mathcal{Y} \Rightarrow [0, 1], \quad P_{Y|X=x}(Y) = P_{Y|X}(Y|X)$$

conditional pmf:

$$P_{Y|X}(Y|X) \stackrel{\text{def}}{=} \frac{P(Y, X)}{P_X(X)} \quad \text{implies} \quad \sum_{Y \in \mathcal{Y}} P_{Y|X}(Y|X) = 1$$

Distribution:

$$P_{Y|X}(Y \in \tilde{E}_Y | X=x) \stackrel{\text{def}}{=} \sum_{Y \in \tilde{E}_Y} P_{Y|X}(Y|X) \quad \text{where } \tilde{E}_Y \subset \mathcal{Y}$$

Continuous  $Y$  for any  $x \in \mathcal{X}$  that  $P_X(x) \neq 0$

$$P_{Y|X=x}: \mathcal{Y} \Rightarrow [0, \infty), \quad P_{Y|X=x}(Y) = P_{Y|X}(Y|X)$$

conditional pdf:

$$P_{Y|X}(y|x) \stackrel{\text{def}}{=} \frac{p(y,x)}{p_X(x)}$$

$$\text{implies } \int_Y P_{Y|X}(y|x) dy = 1$$

Distribution:

$$P_{Y|X}(Y \in E | X=x) \stackrel{\text{def}}{=} \int_E P_{Y|X}(y|x) dy \quad \text{where } E_Y \subset Y$$

Product Rule:

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

More generally:

$$p(x_1, x_2, \dots, x_d) = p(x_d | x_1, \dots, x_{d-1}) \dots p(x_3 | x_1, x_2) p(x_2 | x_1) p(x_1)$$

Bayes' Rule:

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

Note: Sometimes the subscripts are not used for marginal and conditional distributions when it is clear from the context

$$p(x,y) = p_{x,y}(x,y)$$

$$p(x) = p_X(x)$$

$$p(y|x) = p_{Y|X}(y|x)$$

$E_x$ : Probability that I have arthritis given I am young

$$P(Y=1 | X=0) = \sum_{Y \in E_1} P_{Y|X}(y|0)$$

Arthritis Young

|                          |   |
|--------------------------|---|
| $\frac{1}{2}$<br>$(0,0)$ | $\frac{1}{100} (1,0)$<br>$\frac{39}{100} (1,1)$ |
|--------------------------|---|

condition on being young

|                          |                     |
|--------------------------|---------------------|
| $Y=0$<br>$\frac{50}{51}$ | $\frac{1}{51}, Y=1$ |
|--------------------------|---------------------|

$$\begin{aligned}
 &= P_{Y|X}(1|0) \\
 &= \frac{p(0,1)}{P_X(0)} \\
 &= \frac{p(0,1)}{p(0,1) + p(0,0)} \\
 &= \frac{1/100}{1/100 + 1/2} \\
 &= \frac{1/100}{51/100} = \frac{1}{51}
 \end{aligned}$$

Ex: Probability of being young given I have arthritis

$$P(X=0 | Y=1) = P_{X|Y}(0|1)$$

Bayes' Rule  $\Downarrow$

$$\begin{aligned}
 &= \frac{P_{Y|X}(1|0) p(0)}{P_Y(1)} \\
 &= \frac{\frac{1}{51} \frac{51}{100}}{40/100} = \frac{1}{40}
 \end{aligned}$$

Independence: Changing the value of one r.v. doesn't affect the probability of another r.v.

r.v.  $X, Y$  are independent if:  $p(x, y) = p_x(x) p_y(y)$

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

independence implies:  $p(x|y) = p(x)$ ,  $p(y|x) = p(y)$

More generally:

$X_1, X_2, \dots, X_d$  are independent if:  $p(x_1, \dots, x_d) = p(x_1) \dots p(x_d)$

Similarly for distributions:

r.v.  $X, Y$  are independent if:  $P(X \in E_x, Y \in E_y) = P(X \in E_x)P(Y \in E_y)$

$E_x$ :  $X, Y$  are not independent for Arthritis ex

$$p(0, 1) = \frac{1}{100} \neq p_x(0)p_y(1) = \frac{51}{100} \frac{40}{100} = 0.204$$

$E_x$ :  $X_1, X_2 \in \{0, 1\}$  are flips of two different fair coins

$$p(x_1, x_2) = \frac{1}{4} \text{ for all } x_1 \in \mathcal{X}_1, x_2 \in \mathcal{X}_2$$

$$p_{x_1}(x_1)p_{x_2}(x_2) = \frac{1}{2} \frac{1}{2} = \frac{1}{4}$$

|       |   |               |               |
|-------|---|---------------|---------------|
|       |   | $x_2$         |               |
|       |   | H             | T             |
| $x_1$ | H | $\frac{1}{4}$ | $\frac{1}{4}$ |
|       | T | $\frac{1}{4}$ | $\frac{1}{4}$ |

What happens when  $Z=(X,Y)$  with  $Y$  discrete and  $X$  continuous?

$p: \mathcal{X} \times \mathcal{Y} \rightarrow ?$  pmf or pdf? **Ans: neither**

Instead we will write  $p(x,y)$  in terms of a marginal pdf for  $X$  and a conditional pmf for  $Y|X$

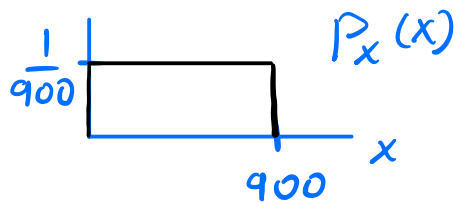
$$p(x,y) = p_X(x) p_{Y|X}(y|x) \quad \text{product rule}$$

where  $p_{Y|X=x}: \mathcal{Y} \rightarrow [0,1]$  is a pmf

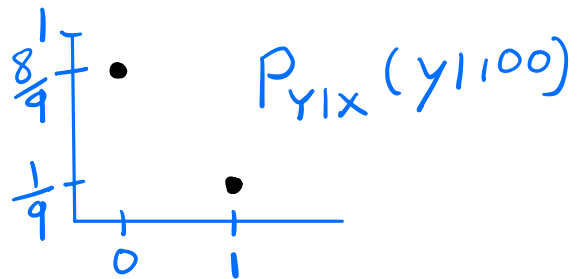
$p_X(x): \mathcal{X} \rightarrow [0,\infty)$  is a pdf

Ex:  $X \in \mathcal{X} = [0, 900]$ ,  $Y \in \mathcal{Y} = \{0, 1\}$  Barolo

pdf:  $p_X = \text{Uniform}(0, 900)$   
 $= \frac{1}{900}$



$p_{Y|X=x} = \text{Bernoulli}\left(\frac{x}{900}\right)$



pmf:  $p_{Y|X}(y|x) = \begin{cases} \frac{x}{900} & \text{if } y=1 \\ 1 - \frac{x}{900} & \text{if } y=0 \end{cases}$   
 Defn of Bernoulli( $\frac{x}{900}$ )

$$\begin{aligned} P(X \in [0, 50], Y=1) &= \int_0^{50} \left( \sum_{Y \in \{1\}} p(y,x) \right) dx \\ &= \int_0^{50} \left( \sum_{Y \in \{1\}} p_{Y|X}(y|x) p_X(x) \right) dx \end{aligned}$$

$$= \int_0^{50} p_{YX}(1|x) p_X(x) dx$$

$$= \int_0^{50} \frac{x}{900} \frac{1}{900} dx$$

$$= \frac{1}{810000} \left. \frac{x^2}{2} \right|_0^{50}$$

$$= \frac{1}{810000} \frac{2500}{2}$$

$$= \frac{1}{648} \quad 0.154\%$$

# Representing Random Features, Labels, and Datasets

## Random variables:

$$D = (Z_1, Z_2, \dots, Z_n) \in \mathcal{Z}_1 \times \dots \times \mathcal{Z}_n = \mathcal{Z}^n \quad \text{since } \mathcal{Z} = \mathcal{Z}_1 = \dots = \mathcal{Z}_n$$

$$Z_i = (\vec{X}_i, Y_i) \in \mathcal{X} \times \mathcal{Y} = \mathcal{Z} \quad \text{each } Z_i \text{ is a feature-label pair}$$

$$\vec{X}_i = (X_{i,1}, \dots, X_{i,d})^T \in \mathbb{R}^d = \mathcal{X} \quad \vec{X}_i \text{ is a feature vector}$$

## Distributions:

$P_D$ : distribution for  $D$ ,  $P_{Z_i}$ : marginal distribution for  $Z_i$

assumptions:

1.  $(\vec{X}_i, Y_i) = Z_i$  are independent for all  $i \in \{1, \dots, n\}$
  2.  $P_{Z_1} = P_{Z_2} = \dots = P_{Z_n} = P_Z$  all  $Z_i$  have the same distribution
- " $(\vec{X}_i, Y_i)$  are independent and identically distributed (i.i.d.)"

$$\begin{aligned} P_D(Z_1 \in \tilde{E}_1, \dots, Z_n \in \tilde{E}_n) &= P_{Z_1}(Z_1 \in \tilde{E}_1) \cdots P_{Z_n}(Z_n \in \tilde{E}_n) \quad \leftarrow 1. \\ &= P_Z(Z_1 \in \tilde{E}_1) \cdots P_Z(Z_n \in \tilde{E}_n) \quad \leftarrow 2. \end{aligned}$$

## Equivalently:

$$D = ((\vec{X}_1, Y_1), \dots, (\vec{X}_n, Y_n)) \in (\mathcal{X} \times \mathcal{Y})^n$$

where  $(\vec{X}_i, Y_i) \sim P_{\vec{X}, Y}$  are independent for all  $i \in \{1, \dots, n\}$   
↑ "sampled/distributed according to"

$D$  contains  $n$  independent samples of  $(\vec{X}_i, Y_i)$   
"feature-label" pairs all coming from the same  
distribution  $P_{\vec{X}, Y}$

## Functions of Random Variables

A function of a r.v. is a r.v.

Ex: ( $X$  is a fair six-sided dice)

$$X \in \{1, 2, 3, 4, 5, 6\} = \mathcal{X} \quad \text{with} \quad p(x) = \frac{1}{6}$$

$$f(X) = X^2 \in \underbrace{\{1^2, 2^2, 3^2, 4^2, 5^2, 6^2\}}_{\text{outcome space for } f(X)} = \mathcal{Y} \quad \text{is a r.v.}$$

Notice  $f: \mathcal{X} \rightarrow \mathcal{Y}$

Sometimes we give the r.v. a new symbol

$$Y = f(X) = X^2$$

$$P_Y(y) = P_{f(X)}(y) = \frac{1}{6} \quad \text{where } y \in \{1, 2^2, 3^2, 4^2, 5^2, 6^2\}$$

In this case  $P_Y(x^2) = p(x)$  where  $x \in \mathcal{X}$

$$\text{ex: } P_Y(9) = P_Y(3^2) = p(3) = \frac{1}{6}$$



Ex: ( $X$  is the payout from a slot machine)

$X \in [-10, 10]$  with  $p(x) = \frac{1}{20}$ ,  $\mathbb{P} = \text{Uniform}(-10, 10)$

$Y = f(X) = X^2 \in [0, 100] = \mathcal{Y}$

$p_Y(y) = \frac{1}{20\sqrt{y}}$  much more complicated

In general  $p_Y$  is complicated and we will not need to know how to calculate it

The Predictor and Learner are functions of r.v.

Ex: (Predictor)

$\vec{X} = (X_1, X_2)^T \in \mathbb{R}^2 = \mathcal{X}$  with  $\mathbb{P}_{\vec{X}}$

predictor:  $f: \mathcal{X} \rightarrow \mathcal{Y}$  where  $\mathcal{Y} = \mathbb{R}$

$f(\vec{X}) = 3 + 6X_1 + 2.5X_2$  is a r.v. with values in  $\mathcal{Y}$

and has some distribution  $\mathbb{P}_{f(\vec{X})}$

Ex: (Learner)

$D = ((\vec{X}_1, Y_1), \dots, (\vec{X}_n, Y_n)) \in (\mathcal{X} \times \mathcal{Y})^n$  with  $\mathbb{P}_D$

Learner:  $\mathcal{A}: (\mathcal{X} \times \mathcal{Y})^n \rightarrow \{f \mid f: \mathcal{X} \rightarrow \mathcal{Y}\} = \mathcal{F}$

$\mathcal{A}(D) = f$  is a r.v. with values in  $\mathcal{F}$

example and has some distribution  $P_{\mathcal{A}(D)}$

if  $D = ((7, 6), (12, 2.5))$  where  $n=2$ ,  $\mathcal{X}=\mathbb{R}$ ,  $\mathcal{Y}=\mathbb{R}$

then  $f_D$  can be  $f(x) = 2.5 + 6x$

This means we can talk about things like:

- What is the probability the Predictor  $f(\tilde{x})$  outputs some value  $y$
- What is the probability the Learner  $\mathcal{A}(D)$  outputs some predictor  $f$

# Expectation and Variance

Expected Value of a r.v.: average value of the r.v.  
if you sample from its distribution infinitely many times.

The r.v. must take values in  $\mathbb{R}$ .

It is not always the value we expect to see most frequently (that is the mode)

$X \in \mathcal{X}$  is a r.v. with pmf or pdf  $p$

$$\mathbb{E}[X] \stackrel{\text{def}}{=} \begin{cases} \sum_{x \in \mathcal{X}} x p(x) & \text{if } X \text{ is discrete} \\ \int_{\mathcal{X}} x p(x) dx & \text{if } X \text{ is continuous} \end{cases}$$

Ex: (fair six-sided dice)

$X \in \{1, 2, 3, 4, 5, 6\} = \mathcal{X}$  and  $P = \text{Uniform}(n=6)$

thus  $p(x) = \frac{1}{6}$

$$\begin{aligned} \mathbb{E}[X] &= \sum_{x \in \mathcal{X}} x p(x) = 1 \cdot \frac{1}{6} + 2 \cdot \frac{1}{6} + 3 \cdot \frac{1}{6} + 4 \cdot \frac{1}{6} + 5 \cdot \frac{1}{6} + 6 \cdot \frac{1}{6} \\ &= 3.5 \end{aligned}$$

This is not an number you can roll on a dice!

$E_x$ : (Unfair coin)

$X \in \{0, 1\}$  and  $P = \text{Bernoulli}(\alpha)$

thus  $p(1) = \alpha, p(0) = 1 - \alpha$

$$E[X] = \sum_{x \in \mathcal{X}} x p(x) = 0 \cdot (1 - \alpha) + 1 \cdot \alpha = \alpha$$

This is not a result of a coin flip (unless  $\alpha = 1$  or  $\alpha = 0$ )

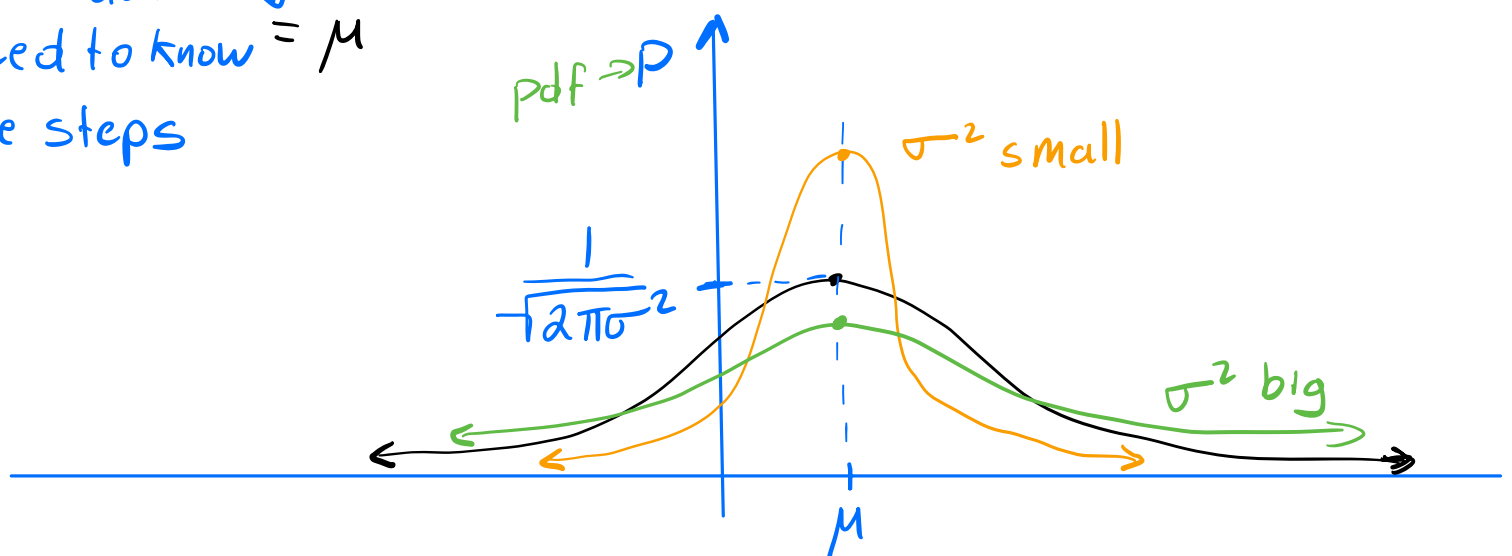
$E_x$ : (Normal distribution)

$X \in \mathbb{R} = \mathcal{X}$  and  $P = \mathcal{N}(\mu, \sigma^2)$

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

$$E[X] = \int_{\mathcal{X}} x \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right) dx$$

You don't need to know the steps  $= \mu$



## Expected value of functions of r.v.:

$X \in \mathcal{X}$  is a r.v. with pmf or pdf  $p$

The function  $f: \mathcal{X} \rightarrow \mathcal{Y}$  must have  $\mathcal{Y} = \mathbb{R}$

$$\mathbb{E}[f(X)] \stackrel{\text{def}}{=} \begin{cases} \sum_{x \in \mathcal{X}} f(x) p(x) & \text{if } X \text{ is discrete} \\ \int_{\mathcal{X}} f(x) p(x) dx & \text{if } X \text{ is continuous} \end{cases}$$

Ex: ( $X$  is the payout from a slot machine)

$X \in [-10, 10]$  with  $p(x) = \frac{1}{20}$ ,  $P = \text{Uniform}(-10, 10)$

$Y = f(X) = X^2 \in [0, 100] = \mathcal{Y}$

$p_Y(y) = \frac{1}{20\sqrt{y}}$  much more complicated

$$\mathbb{E}[f(X)] = \int_{\mathcal{X}} f(x) p(x) dx = \int_{-10}^{10} x^2 \frac{1}{20} dx$$

$$= \frac{x^3}{3} \cdot \frac{1}{20} \Big|_{-10}^{10}$$

$$= \left( \frac{1000}{3} - \frac{(-1000)}{3} \right) \cdot \frac{1}{20}$$

$$= \frac{2000}{60} = 33.333$$

It turns out

$$\mathbb{E}[f(x)] = \mathbb{E}[Y] = \int_{\mathcal{Y}} y P_Y(y) dy$$

exercise  $\rightarrow$   $\approx 33.333$

Usually we don't know  $P_Y = P_{f(x)}$

So we work with  $p$

Variance of a r.v.: How much the r.v. varies from its expected value on average

$X \in \mathcal{X}$  is a r.v. with pmf or pdf  $p$

$$\text{Var}[X] \stackrel{\text{def}}{=} \mathbb{E}\left[\underbrace{(X - \mathbb{E}[X])^2}_{\downarrow}\right] = \mathbb{E}[X^2] - (\mathbb{E}[X])^2$$

this is just a function of the r.v.  $X$

$\mathbb{E}_X$ : (Unfair coin)

$X \in \{0, 1\}$  and  $\mathbb{P} = \text{Bernoulli}(\alpha)$

thus  $p(1) = \alpha, p(0) = 1 - \alpha$

$$\begin{aligned} E[X] &= \sum_{x \in \mathcal{X}} x p(x) = 0 \cdot (1-\alpha) + 1 \cdot \alpha \\ &= \alpha \end{aligned}$$

$$\begin{aligned} \text{Var}[X] &= E[(X - E[X])^2] \\ &= \sum_{x \in \mathcal{X}} (x - E[X])^2 p(x) \\ &= (0 - \alpha)^2 \cdot (1 - \alpha) + (1 - \alpha)^2 \cdot \alpha \\ &= \alpha^2 - \alpha^3 + \alpha - 2\alpha^2 + \alpha^3 \\ &= \alpha - \alpha^2 \\ &= \alpha(1 - \alpha) \end{aligned}$$

or

$$\begin{aligned} \text{Var}[X] &= E[X^2] - (E[X])^2 \\ &= \sum_{x \in \mathcal{X}} x^2 p(x) - \alpha^2 \\ &= 0^2 \cdot (1 - \alpha) + 1^2 \cdot \alpha - \alpha^2 \\ &= \alpha(1 - \alpha) \end{aligned}$$

Ex: (Normal distribution)

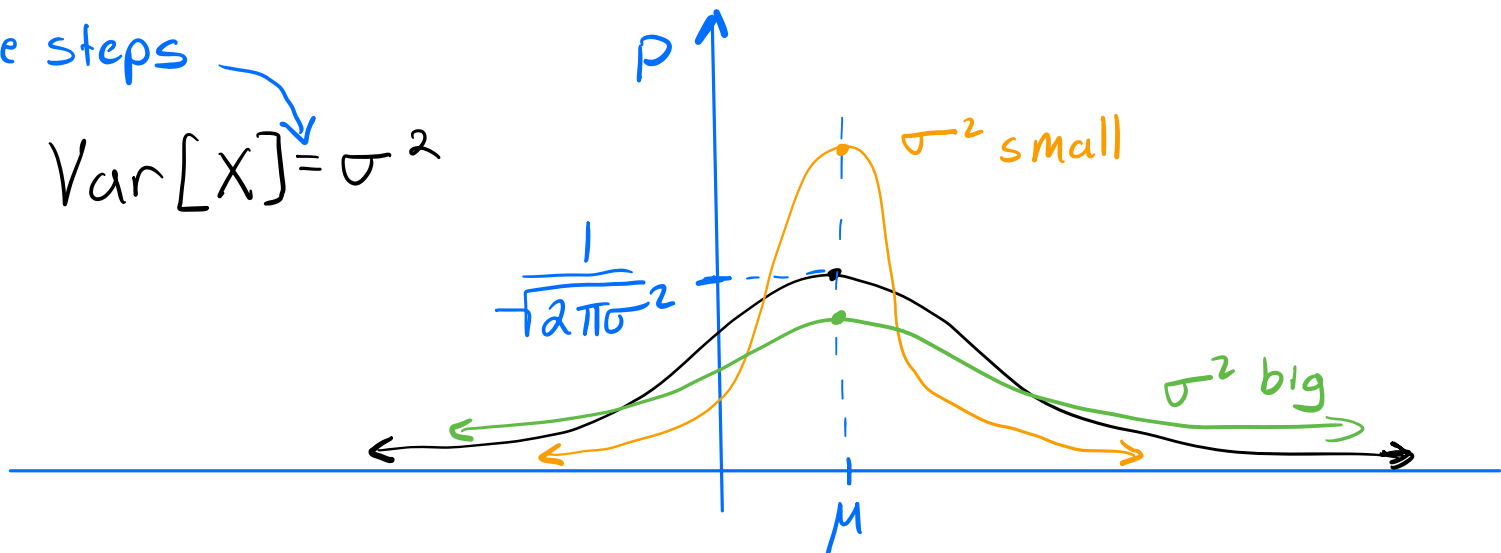
$$X \in \mathbb{R} = \mathcal{X} \text{ and } \mathbb{P} = \mathcal{N}(\mu, \sigma^2)$$

$$\text{thus } p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right)$$

$$E[X] = \int_{-\infty}^{\infty} x \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x-\mu)^2\right) dx$$

You don't need to know the steps  $\rightarrow \mu$

$$\text{Var}[X] = \sigma^2$$



Multivariate Expected Value:

$Z = (X, Y) \in \mathcal{X} \times \mathcal{Y} = \mathcal{Z}$  is a r.v.

$$f: \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$$

$$E[f(X, Y)] = \begin{cases} \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} f(x, y) p(x, y) & \text{if } X, Y \text{ are discrete} \\ \int_{\mathcal{X}} \int_{\mathcal{Y}} f(x, y) p(x, y) dy dx & \text{if } X, Y \text{ are continuous} \\ \int_{\mathcal{X}} \left( \sum_{y \in \mathcal{Y}} f(x, y) p(y|x) \right) p(x) dx & \text{if } Y \text{ is discrete and } X \text{ is continuous} \\ \sum_{x \in \mathcal{X}} \left( \int_{\mathcal{Y}} f(x, y) p(y|x) dy \right) p(x) & \text{if } Y \text{ is continuous and } X \text{ is discrete} \end{cases}$$

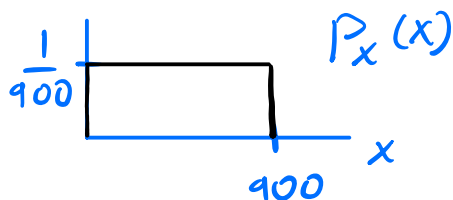


you can always use:

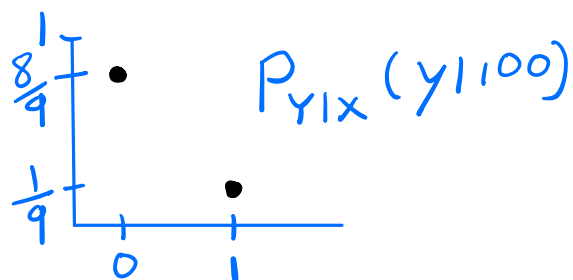
$$p(x, y) = p(y|x)p(x) = p(x|y)p(y)$$

Ex:  $X \in \mathcal{X} = [0, 900]$ ,  $Y \in \mathcal{Y} = \{0, 1\}$  ← Barolo

pdf:  $p_X = \text{Uniform}(0, 900)$   
 $= \frac{1}{900}$



$$P_{Y|X=x} = \text{Bernoulli}\left(\frac{x}{900}\right)$$



pmf:  $P_{Y|X}(y|x) = \begin{cases} \frac{x}{900} & \text{if } y=1 \\ 1 - \frac{x}{900} & \text{if } y=0 \end{cases}$

Defn of Bernoulli( $\frac{x}{900}$ )

$$f(x, Y) = \left(\frac{x}{900} - Y\right)^2$$

$$E[f(x, Y)] = \int_{\mathcal{X}} \left( \sum_{Y \in \mathcal{Y}} f(x, Y) p(Y|x) \right) p(x) dx$$

$$= \int_0^{900} \left( \sum_{Y \in \{0, 1\}} \left(\frac{x}{900} - Y\right)^2 p(Y|x) \right) p(x) dx$$

$$= \int_0^{900} \left( \left(\frac{x}{900} - 0\right)^2 \left(1 - \frac{x}{900}\right) + \left(\frac{x}{900} - 1\right)^2 \left(\frac{x}{900}\right) \right) \frac{1}{900} dx$$

$$= \frac{1}{900} \int_0^{900} \frac{x}{900} \left(1 - \frac{x}{900}\right) dx$$

$$= \frac{1}{900} \left( \frac{x^2}{1800} \Big|_0^{900} - \frac{x^3}{3 \cdot 900^2} \Big|_0^{900} \right)$$

$$= \frac{1}{6}$$

Conditional Expected Value:

$(X, Y) \in \mathcal{X} \times \mathcal{Y}$  is a r.v.

$P = P_{Y|X}$  is a conditional pmf or pdf

$f: \mathcal{Y} \rightarrow \mathbb{R}$

$$\mathbb{E}[f(Y) | X=x] = \begin{cases} \sum_{y \in \mathcal{Y}} f(y) p(y|x) & \text{if } Y \text{ is discrete} \\ \int_{\mathcal{Y}} f(y) p(y|x) & \text{if } Y \text{ is continuous} \end{cases}$$

## Useful Properties

Let  $X, Y$  be r.v. and  $c \in \mathbb{R}$  be a constant

$$1. \mathbb{E}[cX] = c \mathbb{E}[X]$$

$$2. \mathbb{E}[X+Y] = \mathbb{E}[X] + \mathbb{E}[Y]$$

$$3. \mathbb{E}[Y] = \mathbb{E}[\mathbb{E}[Y|X]]$$

$$4. \text{Var}[c] = 0$$

$$5. \text{Var}[cX] = c^2 \text{Var}[X]$$

If  $X$  and  $Y$  are independent:

$$6. \mathbb{E}[XY] = \mathbb{E}[X] \mathbb{E}[Y]$$

$$7. \text{Var}[X+Y] = \text{Var}[X] + \text{Var}[Y]$$