

# 1 Review of probability theory

## 1.1 Why probability theory?

- Probability theory provides mathematical framework for reasoning about prediction problems
- (Some alternatives: approximation theory, game theory, ...)
- Basic idea: regard quantities you are uncertain about (e.g., quantities you want to predict) as random variables defined on a probability space
- Starting from basic idea, can use probability theory to derive properties of optimal predictions, characterize uncertainty of error rate estimates, design and analyze learning algorithms, etc.

## 1.2 Probability spaces

- Goal: mathematical model for experiment with random outcomes (E.g., coin tosses, dice rolls, roulette wheel spins, ...)
- A (discrete) probability space  $(\Omega, m)$  is comprised of a sample space  $\Omega$  and a probability (mass) function  $m$ 
  - Sample space  $\Omega$  is the (finite or countable) set of possible outcomes
  - An event is a subset of  $\Omega$
  - Example: toss a coin
    - \* Possible outcomes:  $\Omega = \{\text{H}, \text{T}\}$
    - \* “heads” =  $\{\text{H}\}$
    - \* “tails” =  $\{\text{T}\}$
    - \* ...
  - Example: toss a coin twice
    - \* Possible outcomes:  $\Omega = \{\text{TT}, \text{TH}, \text{HT}, \text{HH}\}$
    - \* “both tails” =  $\{\text{TT}\}$
    - \* “at least one heads” =  $\{\text{TH}, \text{HT}, \text{HH}\}$
    - \* ...

- Example: roll a 6-sided die
  - \* Possible outcomes:  $\Omega = \{\square, \square, \square, \square, \square, \square\}$
  - \* “odd” =  $\{\square, \square, \square\}$
  - \* “even” =  $\{\square, \square, \square\}$
  - \* “at most 3” =  $\{\square, \square, \square\}$
  - \* ...
- Example: repeatedly roll a 6-sided die and stop after seeing a “6”
  - \* Possible outcomes:  $\Omega = \{\square, \square\square, \square\square, \square\square, \square\square, \square\square, \square\square\square, \square\square\square, \dots\}$
  - \* “one roll” =  $\{\square\}$
  - \* “two rolls” =  $\{\square\square, \square\square, \square\square, \square\square\}$
  - \* ...
- Probability (mass) function  $m$  is a function that assigns a real number to each outcome in  $\Omega$  in a way that satisfies
  - \*  $m(\omega) \geq 0$  for all  $\omega \in \Omega$  (non-negativity), and
  - \*  $\sum_{\omega \in \Omega} m(\omega) = 1$  (normalization)
- Probability of an event  $E \subseteq \Omega$  in probability space  $(\Omega, m)$  is

$$\Pr(E) = \sum_{\omega \in E} m(\omega)$$

(Notation unfortunately does not explicitly show  $(\Omega, m)$ )

- Sometimes we “abuse notation” by writing  $m(E)$  to mean  $\Pr(E)$
- Example: toss a fair coin twice
  - \*  $m(\omega) = 1/4$  for every possible outcome  $\omega$

$$\Pr(\text{tosses come up on same side}) = \frac{1}{4} + \frac{1}{4} = \frac{1}{2}$$

- Some events can be described in terms of other events using set theory

- Union (“or”)

$$A \cup B = \{\omega \in \Omega : \omega \in A \text{ or } \omega \in B\}$$

- Intersection (“and”)

$$A \cap B = \{\omega \in \Omega : \omega \in A \text{ and } \omega \in B\}$$

– Complement (“not”)

$$A^c = \{\omega \in \Omega : \omega \notin A\}$$

– Difference (“and not”)

$$A - B = \{\omega \in \Omega : \omega \in A \text{ and } \omega \notin B\}$$

(sometimes also written “ $A \setminus B$ ”; same as  $A \cap B^c$ )

– Example: roll a fair 6-sided die twice

\*  $A$  = “first roll is even”

\*  $B$  = “second roll is at most 3”

\*  $A^c$  = “first roll is odd”

\*  $B^c$  = “second roll is at least 4”

\* “first roll is even and second roll is at most 3”

$$A \cap B = \{\square\square, \square\square, \square\square, \square\square, \square\square, \square\square, \square\square, \square\square, \square\square\}$$

so

$$\Pr(A \cap B) = \frac{9}{36} = \frac{1}{4}$$

\* “first toss is even or second toss is at most 3”

$$\begin{aligned} A \cup B &= (A^c \cap B^c)^c \\ &= \Omega - (A^c \cap B^c) \\ &= \Omega - \{\square\square, \square\square, \square\square, \square\square, \square\square, \square\square, \square\square, \square\square, \square\square\} \end{aligned}$$

so

$$\Pr(A \cup B) = \frac{36 - 9}{36} = \frac{3}{4}$$

Q. Suppose a 6-sided die is weighted so that, for each  $k \in \{1, 2, 3, 4, 5, 6\}$ , the side showing  $k$  pips is  $k$  times as likely to show up as the side showing 1 pip. What is the probability that a roll of this die shows an even number of pips?

Q. Suppose  $A$  and  $B$  are events from a probability space such that  $\Pr(A \cap B) = 1/4$ ,  $\Pr(A^c) = 1/3$ , and  $\Pr(B) = 1/2$ . What is  $\Pr(A \cup B)$ ?

### 1.3 Conditional probability

- Suppose, in an experiment described by probability space  $(\Omega, m)$ , you learn that an event  $E$  has occurred, but nothing else
  - Exact outcome  $\omega$  may not yet be known to you
  - What probability space now models the experiment in light of the new information?
  - Conditioning on  $E$ : incorporating information that  $E$  occurred

- New probability space  $(\Omega, m_E)$  with probability function defined in terms of original:

$$m_E(\omega) = \begin{cases} \frac{m(\omega)}{m(E)} & \text{if } \omega \in E \\ 0 & \text{if } \omega \notin E \end{cases}$$

(We require  $m(E) > 0$  in order to define  $m_E$ )

- Can check that  $m_E$  is a valid probability function on  $\Omega$   
(Normalization is ensured by the division by  $\Pr(E)$ )
- Notation: write  $\Pr(F | E)$  for probability of event  $F$  in probability space  $(\Omega, m_E)$ , a.k.a. probability of  $F$  conditioned on  $E$ , a.k.a. (conditional) probability of  $F$  given  $E$
- Example: roll a fair 6-sided die
  - $E = \{\square, \boxtimes, \boxplus\}$  = “even”,  $F = \{\square, \boxtimes, \boxtimes\}$  = “prime”
  - Suppose you learn  $E$  occurred
    - \* Given this information, what is probability of  $F$ ?

$$\Pr(F | E) = \sum_{\omega \in F} m_E(\omega) = \sum_{\omega \in F \cap E} \frac{m(\omega)}{\Pr(E)} = \frac{1/6}{1/2} = \frac{1}{3}$$

- Useful formula for conditional probability:

$$\Pr(F | E) \Pr(E) = \Pr(F \cap E)$$

- Bayes' rule: relates probabilities of event  $F$  before and after conditioning on information that event  $E$  occurs

$$\Pr(F | E) = \Pr(F) \times \frac{\Pr(E | F)}{\Pr(E)}$$

- $\Pr(F | E)$  is probability of  $F$  after conditioning on information that  $E$  occurred
- $\Pr(F)$  is probability of  $F$  in original probability space (before observing that  $E$  occurred)
- Ratio  $\Pr(E | F)/\Pr(E)$  is what relates these probabilities
  - \* Always non-negative, but can be zero (even if  $\Pr(E) > 0$ )
  - \* Whether it is more or less than 1 determines whether probability of  $F$  increases or decreases after incorporating information that  $E$  occurred
- Example: A casino has 100 identically-looking slot machines; but unbeknownst to you, the first 75 are “fair”, and rest are “rigged”. If you play on a “fair” machine, you are equally likely to win or lose. If you play on a “rigged” machine, you always lose.

Suppose you enter the casino, pick a slot machine uniformly at random, play it once, and lose. Given this information, what is the probability that you played on a “rigged” machine?

- Sample space:  $\Omega = \{1, 2, \dots, 100\} \times \{\text{win}, \text{lose}\}$   
(Other choices could also work)
- Events of interest:
  - \*  $R = \{(a, b) \in \Omega : 76 \leq a \leq 100\}$
  - \*  $L = \{(a, b) \in \Omega : b = \text{lose}\}$
- Probabilities of interest:
  - \*  $\Pr(R) = 25/100, \Pr(R^c) = 75/100$
  - \*  $\Pr(L | R) = 1, \Pr(L | R^c) = 1/2$

\* We also need  $\Pr(L)$ :

$$\begin{aligned}\Pr(L) &= \Pr(L \cap R) + \Pr(L \cap R^c) \\ &= \Pr(L | R) \times \Pr(R) + \Pr(L | R^c) \times \Pr(R^c) \\ &= 1 \times \frac{25}{100} + \frac{1}{2} \times \frac{75}{100} \\ &= \frac{1}{4} + \frac{3}{8} = \frac{5}{8}\end{aligned}$$

– Using Bayes' rule:

$$\begin{aligned}\Pr(R | L) &= \Pr(R) \times \frac{\Pr(L | R)}{\Pr(L)} \\ &= \frac{1}{4} \times \frac{1}{5/8} \\ &= \frac{2}{5} = 40\%\end{aligned}$$

– Before playing the machine, 25% probability that picked machine is rigged; after playing machine and observing that you lost, the probability has increased to 40%

Q. In the casino example, suppose you play a randomly picked machine two times, and lose both times. What is probability that you picked a rigged machine, given this information?

Q. You repeatedly roll a fair 6-sided die and stop after seeing 6 pips face up. Suppose only even numbers of pips show up in all rolls. What is the probability that the number of rolls is 1, given this information?

## 1.4 Random variables

- Random variable  $X$  (on  $(\Omega, m)$ ) is a function that assigns a real number to each outcome in  $\Omega$

- Facilitates quantitative analysis of experiments modeled by probability spaces

- $X$  defines probability space  $(\mathbb{R}, p_X)$  with  $p_X$  defined by

$$p_X(x) = \Pr(X = x) = \Pr(\{\omega \in \Omega : X(\omega) = x\})$$

- \*  $p_X$  is the probability (mass) function for  $X$
- \* Say  $p_X$  specifies the probability distribution of  $X$
- \* Sample space is set of real numbers  $\mathbb{R}$ , but probability function  $p_X$  takes values 0 outside of range of  $X$

$$\text{range}(X) = \{X(\omega) : \omega \in \Omega\}$$

- \* So can also regard sample space as  $\text{range}(X)$
  - \* Shorthand: “ $X = x$ ” means  $\{\omega \in \Omega : X(\omega) = x\}$
- Example: toss a fair coin three times
- \*  $X$  = number of heads

$$\begin{array}{ll} X(\text{TTT}) = 0 & X(\text{TTH}) = 1 \\ X(\text{THT}) = 1 & X(\text{THT}) = 2 \\ X(\text{HTT}) = 1 & X(\text{HTH}) = 2 \\ X(\text{HHT}) = 2 & X(\text{HHH}) = 3 \end{array}$$

$$\begin{array}{c|cccc} x & 0 & 1 & 2 & 3 \\ \hline p_X(x) & 1/8 & 3/8 & 3/8 & 1/8 \end{array}$$

- \* Event “at least one heads” is also written as “ $X \geq 1$ ”

$$\Pr(X \geq 1) = \frac{7}{8}$$

- Example: roll a fair 6-sided die twice
- \*  $X$  = number of pips from the first roll
  - \*  $Y$  = number of pips from the second roll
  - \*  $Z = X + Y$

$$Z(\begin{array}{|c|c|} \hline \square & \square \\ \hline \end{array}) = X(\begin{array}{|c|c|} \hline \square & \square \\ \hline \end{array}) + Y(\begin{array}{|c|c|} \hline \square & \square \\ \hline \end{array}) = 1 + 1 = 2$$

$$Z(\begin{array}{|c|c|} \hline \square & \blacksquare \\ \hline \end{array}) = X(\begin{array}{|c|c|} \hline \square & \blacksquare \\ \hline \end{array}) + Y(\begin{array}{|c|c|} \hline \square & \blacksquare \\ \hline \end{array}) = 1 + 2 = 3$$

$$Z(\begin{array}{|c|c|} \hline \square & \blacksquare \blacksquare \\ \hline \end{array}) = X(\begin{array}{|c|c|} \hline \square & \blacksquare \blacksquare \\ \hline \end{array}) + Y(\begin{array}{|c|c|} \hline \square & \blacksquare \blacksquare \\ \hline \end{array}) = 1 + 3 = 4$$

etc.

- Expectation (a.k.a. expected value, mean, average) of random variable  $X$  in probability space  $(\Omega, m)$

$$\mathbb{E}(X) = \sum_{\omega \in \Omega} X(\omega)m(\omega)$$

- Often more convenient to use equivalent formula

$$\mathbb{E}(X) = \sum_x x \Pr(X = x) = \sum_x x p_X(x)$$

(Summation is taken over  $x \in \text{range}(X)$ )

- Example: roll a fair 6-sided die

- \*  $X$  = number of pips

$$\begin{aligned} \mathbb{E}(X) &= \sum_{x=1}^6 x p_X(x) \\ &= 1 \times \frac{1}{6} + 2 \times \frac{1}{6} + 3 \times \frac{1}{6} + 4 \times \frac{1}{6} + 5 \times \frac{1}{6} + 6 \times \frac{1}{6} \\ &= \frac{21}{6} = 3.5 \end{aligned}$$

- Example: toss a fair coin three times

- \*  $X$  = number of heads

$$\begin{aligned} \mathbb{E}(X) &= \sum_{x=0}^3 x p_X(x) \\ &= 0 \times \frac{1}{8} + 1 \times \frac{3}{8} + 2 \times \frac{3}{8} + 3 \times \frac{1}{8} \\ &= \frac{12}{8} = 1.5 \end{aligned}$$

- If  $X$  is a random variable, and  $Y = aX + b$  for some real numbers  $a$  and  $b$ , then

$$\mathbb{E}(Y) = \mathbb{E}(aX + b) = a \mathbb{E}(X) + b$$

- Caution: not all random variables have an expectation

- \* Example:  $p_X(x) = 1/x - 1/(x+1)$  for all positive integers  $x$

- Beyond the expected value

- Random variables with same expected value can be very different
  - Example:

- \* Toss fair coin 5 times;  $X =$  number of heads,  $\mathbb{E}(X) = 2.5$ 
  - $\text{range}(X) = \{0, 1, 2, 3, 4, 5\}$
  - $|\{\omega \in \Omega : X(\omega) \in \{2, 3\}\}| = 20$
  - $|\{\omega \in \Omega : X(\omega) \in \{0, 1, 4, 5\}\}| = 12$
  - So values “close” to the expectation are more likely than those “far” from the expectation
- \* Roll a fair 6-sided die;  $Y =$  number of pips  $- 1$ ,  $\mathbb{E}(Y) = 2.5$ 
  - All possible values  $\{0, 1, 2, 3, 4, 5\}$  of  $Y$  are equally likely, regardless of distance to the expectation
- \*  $X$  is less “spread out” than  $Y$
- Variance: convenient measure of a random variable’s “spread”

$$\text{var}(X) = \mathbb{E}((X - \mu)^2)$$

where  $\mu = \mathbb{E}(X)$

- \* The square-root of  $\text{var}(X)$ —called standard deviation—is roughly how much  $X$  deviates from  $\mu$  on average
  - $\text{stddev}(X) = \sqrt{\text{var}(X)}$
  - Caveat:  $\sqrt{\mathbb{E}(X^2)}$  is not necessarily the same as  $\mathbb{E}(\sqrt{X^2})$
- \*  $\mathbb{E}(|X - \mu|)$  is exactly how much  $X$  deviates from  $\mu$  on average, but less convenient to work with mathematically
- If  $X$  is a random variable, and  $Y = aX + b$  for some real numbers  $a$  and  $b$ , then

$$\text{var}(Y) = \text{var}(aX + b) = \text{var}(aX) = a^2 \text{var}(X)$$

- There are many other “summary statistics” for random variables

- Q. You repeatedly roll a fair 6-sided die and stop after seeing 6 pips face up. What is the expected number of rolls?
- Q. If  $X$  is the number of heads in 5 tosses of a fair coin, and  $Y$  is number of pips shown in the roll of a fair 6-sided die, what are the variances of  $X$  and  $Y$ ?

## 1.5 Multiple random variables

- If each of  $X$  and  $Y$  is a random variable (on  $(\Omega, m)$ ), then (2-dimensional) random vector  $Z = (X, Y)$  is a  $\mathbb{R}^2$ -valued function on  $\Omega$  given by

$$Z(\omega) = (X(\omega), Y(\omega))$$

- $Z$  defines probability space  $(\mathbb{R}^2, p_Z)$  by

$$p_Z(x, y) = \Pr(X = x \wedge Y = y)$$

$p_Z$  is joint probability function for  $(X, Y)$

- Example: roll a fair 6-sided die 10 times
  - \*  $X$  = number of rolls with 6 pips
  - \*  $Y$  = number of rolls with 5 pips
- Can generalize to  $n$ -tuples of random variables to get  $n$ -dimensional random vectors
- Useful fact: If  $X$  and  $Y$  are random variables (on  $(\Omega, m)$ ), then

$$\mathbb{E}(X + Y) = \mathbb{E}(X) + \mathbb{E}(Y)$$

i.e., expectation is additive

- Example: roll a fair 6-sided die 10 times
  - \*  $X$  = number of rolls with 6 pips
  - \*  $Y$  = number of rolls with 5 pips
  - \*  $\mathbb{E}(X + Y) = \mathbb{E}(X) + \mathbb{E}(Y) = 5/3 + 5/3 = 10/3$
- Generalizes to sums of  $n$  random variables  $X_1, \dots, X_n$

$$\mathbb{E}(X_1 + \dots + X_n) = \mathbb{E}(X_1) + \dots + \mathbb{E}(X_n)$$

and also linear combinations

$$\mathbb{E}(a_1 X_1 + \dots + a_n X_n) = a_1 \mathbb{E}(X_1) + \dots + a_n \mathbb{E}(X_n)$$

i.e., expectation is linear

- Random variables  $X$  and  $Y$  are independent if, for all pairs of real numbers  $(x, y)$ ,

$$\Pr(X = x \wedge Y = y) = \Pr(X = x) \times \Pr(Y = y)$$

i.e.,

$$p_{(X,Y)}(x, y) = p_X(x)p_Y(y) \quad \text{for all } (x, y)$$

– Example: roll a fair 6-sided die

$$* X = \begin{cases} 1 & \text{if number of pips is at most 4} \\ 0 & \text{otherwise} \end{cases}$$

$$p_X(0) = \frac{1}{3}, \quad p_X(1) = \frac{2}{3}$$

$$* Y = \begin{cases} 1 & \text{if number of pips is even} \\ 0 & \text{otherwise} \end{cases}$$

$$p_Y(0) = p_Y(1) = \frac{1}{2}$$

\* Joint probability function

$$\begin{array}{c|cccc} (x, y) & (0, 0) & (0, 1) & (1, 0) & (1, 1) \\ \hline p_{(X,Y)}(x, y) & 1/6 & 1/6 & 1/3 & 1/3 \end{array}$$

\* Check that this satisfies  $p_{(X,Y)}(x, y) = p_X(x)p_Y(y)$  for all  $(x, y)$

\* So  $X$  and  $Y$  are independent

\* Note: Here,  $X$  and  $Y$  are special kinds of random variables called indicator random variables—each one indicates whether or not a particular event occurs

\* Notation:

$$X = \mathbb{1}\{\text{number of pips is at most 4}\}$$

$$Y = \mathbb{1}\{\text{number of pips is even}\}$$

\* Distribution of an indicator random variable  $X$  is Bernoulli, written as  $X \sim \text{Bernoulli}(\theta)$ , where  $\theta = \Pr(X = 1)$

– A non-example: roll a fair 6-sided die 10 times

- \*  $X$  = number of rolls with 6 pips
- \*  $Y$  = number of rolls with 5 pips
- \*  $\Pr(X = 10 \wedge Y = 10) = 0$ , yet

$$\Pr(X = 10) = \Pr(Y = 10) > 0$$

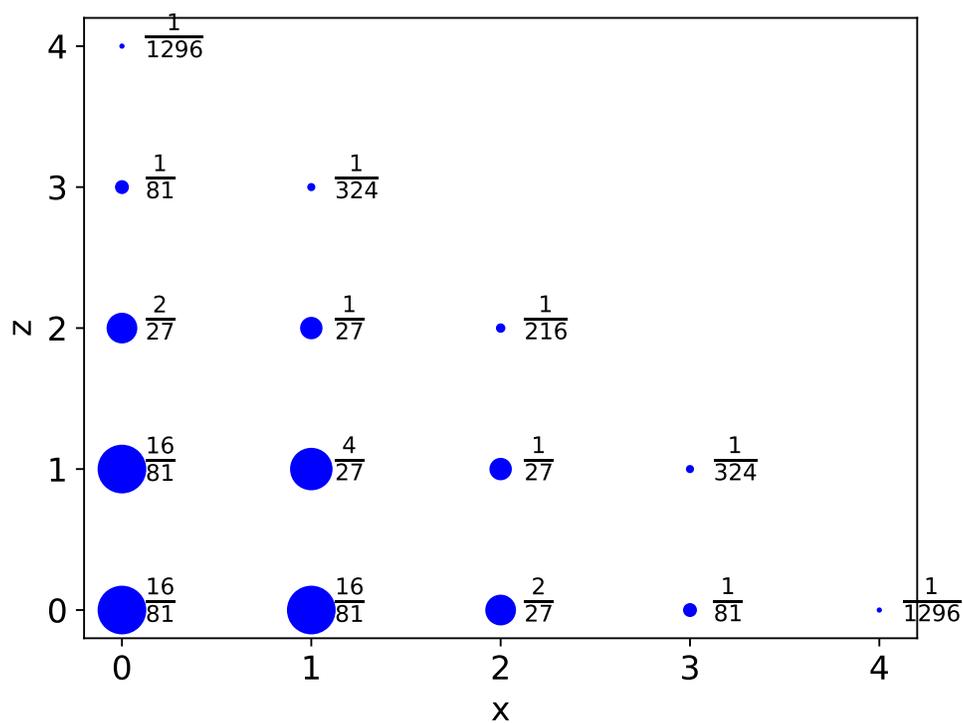
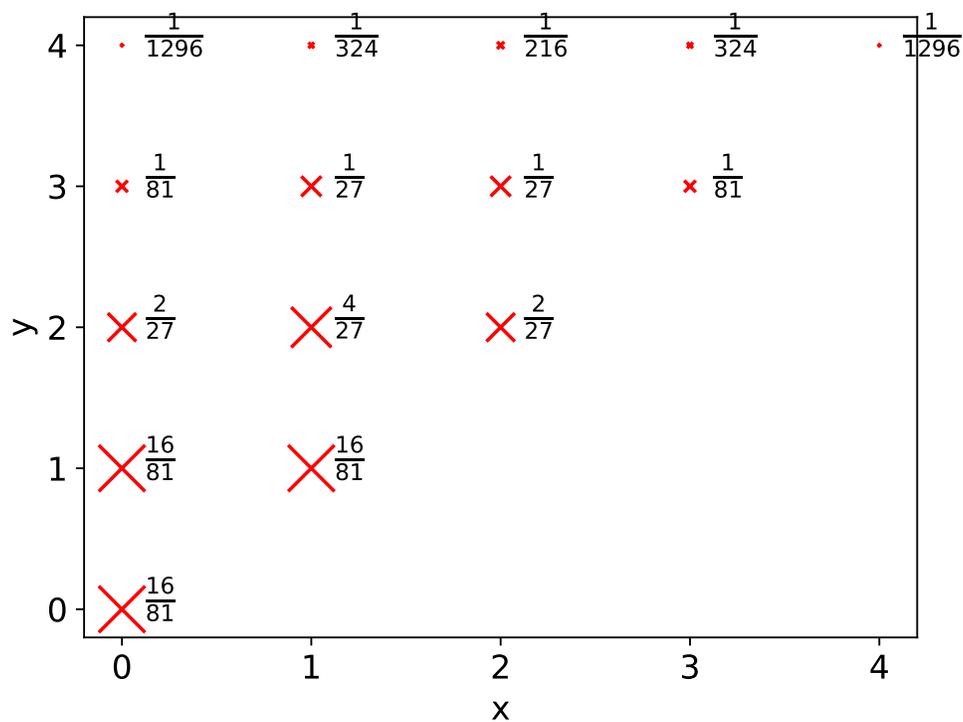
- \* So  $X$  and  $Y$  are not independent
- Generalizes to  $n$  random variables:  $X_1, \dots, X_n$  are independent if, for all  $n$ -tuples of real numbers  $(x_1, \dots, x_n)$ ,

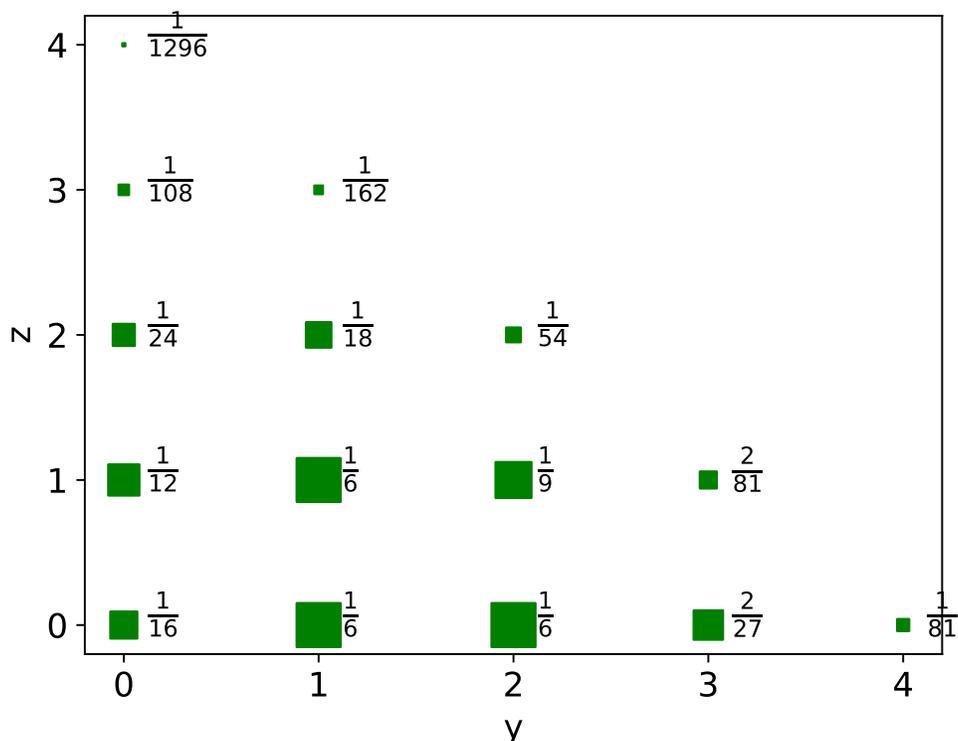
$$\Pr(X_1 = x_1 \wedge \dots \wedge X_n = x_n) = \Pr(X_1 = x_1) \times \dots \times \Pr(X_n = x_n)$$

- Q. Roll a fair 6-sided die; let  $X$  indicate if number of pips is at most 4, and let  $Y$  indicate if number of pips is even. Are  $X$  and  $Y$  independent?
- Q. Toss a fair coin 10 times, and let  $X$  be the number times **HTH** appears as a substring of the outcome. What is the expected value of  $X$ ? (Hint: Write  $X$  as a sum of 8 indicator random variables, and use the linearity of expectation.)

## 1.6 Dependence

- Random variables that are not independent are said to be dependent
- Many different “types” of dependence
  - Example: Roll a fair 6-sided die  $n$  times; let  $X$  be the number of times a  comes up; let  $Y$  be the number of times a  or  comes up; let  $Z$  be the number of times a  comes up
    - \* The larger  $X$  is, the larger  $Y$  must be
    - \* But the larger  $X$  or  $Y$  is, the smaller  $Z$  must be





- Say  $X$  and  $Y$  are positively correlated if  $\mathbb{E}(XY) > \mathbb{E}(X)\mathbb{E}(Y)$
- In die rolling example with  $n = 4$ :

$$\mathbb{E}(XY) = 4/3$$

$$\mathbb{E}(X) = 2/3$$

$$\mathbb{E}(Y) = 4/3$$

$$\mathbb{E}(X)\mathbb{E}(Y) = 8/9$$

So  $X$  and  $Y$  are positively correlated

- Say  $X$  and  $Y$  are negatively correlated if  $\mathbb{E}(XY) < \mathbb{E}(X)\mathbb{E}(Y)$
- Say  $X$  and  $Y$  are uncorrelated if  $\mathbb{E}(XY) = \mathbb{E}(X)\mathbb{E}(Y)$
- If  $X$  and  $Y$  are independent, then they are uncorrelated
- But converse is not necessarily true
- Example: toss a fair coin two times

\*  $X$  = number of heads

$$* Y = \begin{cases} 1 & \text{if first toss is heads and second toss is tails} \\ 0 & \text{if both tosses are the same} \\ -1 & \text{if first toss is tails and second toss is heads} \end{cases}$$

- \*  $\mathbb{E}(X) = 1, \mathbb{E}(Y) = 0, \mathbb{E}(XY) = 0$
- \* So  $X$  and  $Y$  are uncorrelated, but

$$\frac{1}{4} = \Pr(X = 0, Y = 0) \neq \Pr(X = 0) \times \Pr(Y = 0) = \frac{1}{4} \times \frac{1}{2} = \frac{1}{8}$$

- Also many different ways to “measure” dependence

- Covariance between  $X$  and  $Y$  is

$$\begin{aligned} \text{cov}(X, Y) &= \mathbb{E}((X - \mathbb{E}(X))(Y - \mathbb{E}(Y))) \\ &= \mathbb{E}(XY) - \mathbb{E}(X)\mathbb{E}(Y) \end{aligned}$$

- For any constants  $a, b, c, d$ ,

$$\text{cov}(aX + b, cY + d) = ac \text{cov}(X, Y)$$

- (Pearson’s) correlation between  $X$  and  $Y$  is

$$\text{cor}(X, Y) = \frac{\text{cov}(X, Y)}{\text{stddev}(X) \text{stddev}(Y)}$$

- In die rolling example with  $n = 4$ :

$$\begin{aligned} \text{cov}(X, Z) &= -1/9, \quad \text{var}(X) = \text{var}(Z) = 5/9 \\ \text{cor}(X, Z) &= -1/5 \end{aligned}$$

Q. If  $\text{var}(X + Y) = \text{var}(X) + \text{var}(Y)$ , then what can you say about  $\text{cov}(X, Y)$ ?

Q. If  $X = Y$ , then what is the value of  $\text{cor}(X, Y)$ ?

Q. Is it possible to have  $\text{cor}(X, Y) > 1$ ? What about  $\text{cor}(X, Y) < -1$ ?

## 1.7 Marginal and conditional distributions

- Consider random variables  $X$  and  $Y$  (on  $(\Omega, m)$ )
- Marginal distribution of  $Y$  is the probability distribution given by

$$p_Y(y) = \Pr(Y = y) = \Pr(\{\omega \in \Omega : Y(\omega) = y\})$$

– Law of total probability:

$$p_Y(y) = \Pr(Y = y) = \sum_x \Pr(X = x \wedge Y = y) = \sum_x p_{(X,Y)}(x, y)$$

This process of summing  $p_{(X,Y)}(x, y)$  over all possible values of  $X$  is called marginalization

- Conditional distribution of  $Y$  given  $X = x$  is probability distribution  $p_{Y|X=x}$  given by

$$\begin{aligned} p_{Y|X=x}(y) &= \Pr(Y = y \mid X = x) \\ &= \frac{\Pr(Y = y \wedge X = x)}{\Pr(X = x)} \end{aligned}$$

- Conditional expectation of  $Y$  given  $X = x$  is

$$\mathbb{E}(Y \mid X = x) = \sum_y y p_{Y|X=x}(y)$$

– Example: roll a fair 6-sided die

- \*  $X = \mathbb{1}\{\text{number of pips is more than 4}\}$
- \*  $Y = \text{number of pips}$
- \*  $\mathbb{E}(Y \mid X = 0) = 2.5$
- \*  $\mathbb{E}(Y \mid X = 1) = 5.5$
- \*  $Y' = \mathbb{1}\{\text{number of pips is even}\}$
- \*  $\mathbb{E}(Y' \mid X = 0) = 1/2$
- \*  $\mathbb{E}(Y' \mid X = 1) = 1/2$
- \*  $X' = \mathbb{1}\{\text{number of pips is more than 3}\}$
- \*  $\mathbb{E}(Y' \mid X' = 0) = 1/3$
- \*  $\mathbb{E}(Y' \mid X' = 1) = 2/3$

- Regard  $Z = \mathbb{E}(Y \mid X)$  as a random variable in probability space  $(\mathbb{R}, p_X)$

–  $Z(x) = \mathbb{E}(Y \mid X = x)$

- Expected value of  $\mathbb{E}(Y | X)$  is

$$\begin{aligned}\mathbb{E}(\mathbb{E}(Y | X)) &= \sum_x \mathbb{E}(Y | X = x) p_X(x) \\ &= \sum_x \sum_y y p_{Y|X=x} p_X(x) \\ &= \sum_x \sum_y y p_{(X,Y)}(x, y) \\ &= \sum_y y p_Y(y) \\ &= \mathbb{E}(Y)\end{aligned}$$

This fact is called the tower property of conditional expectation

Q. Toss a fair coin two times; let

- $X$  = number of heads
- $Y = \begin{cases} 1 & \text{if first toss is heads and second toss is tails} \\ 0 & \text{if both tosses are the same} \\ -1 & \text{if first toss is tails and second toss is heads} \end{cases}$

For each  $x \in \text{range}(X)$ , what is the expected value of  $Y$  given  $X = x$ ?

## 1.8 Continuous random variables

- So far, we have only considered discrete random variables (which have finite or countable ranges)
  - Probability distribution of random variable  $X$  can be specified either by its probability mass function  $p_X$  or by its (cumulative) distribution function (cdf)  $\text{cdf}_X$ 

$$\text{cdf}_X(x) = \Pr(X \leq x)$$
- A random variable is continuous if its distribution function is a continuous function
  - In some cases, these arise by starting with discrete distributions and taking an appropriate limit

- In this class, we'll only discuss continuous random variables  $X$  whose distribution functions can be written as

$$\text{cdf}_X(x) = \int_{-\infty}^x p_X(u) \, du$$

for a function  $p_X$  called the (probability) density function (pdf)

- Important example: uniform (on unit interval) random variable

$$p_X(x) = \begin{cases} 1 & \text{if } 0 \leq x \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

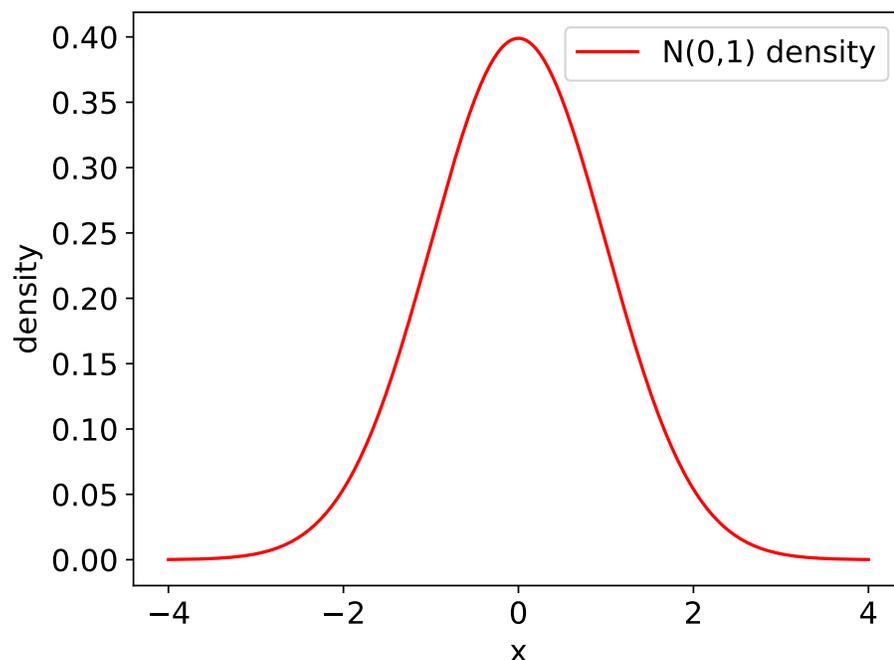
- Notation:  $X \sim \text{Unif}([0, 1])$
- For any subinterval  $I \subseteq [0, 1]$ ,  $\Pr(X \in I)$  is the length of the interval

Uniform (on unit square) random vector:

$$p_{(X,Y)}(x, y) = \begin{cases} 1 & \text{if } 0 \leq x \leq 1 \text{ and } 0 \leq y \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

- Notation:  $(X, Y) \sim \text{Unif}([0, 1]^2)$
- Can verify that  $X$  and  $Y$  are independent, and each of  $X$  and  $Y$  has marginal distribution  $\text{Unif}([0, 1])$
- Another important example: a standard normal random variable has density function

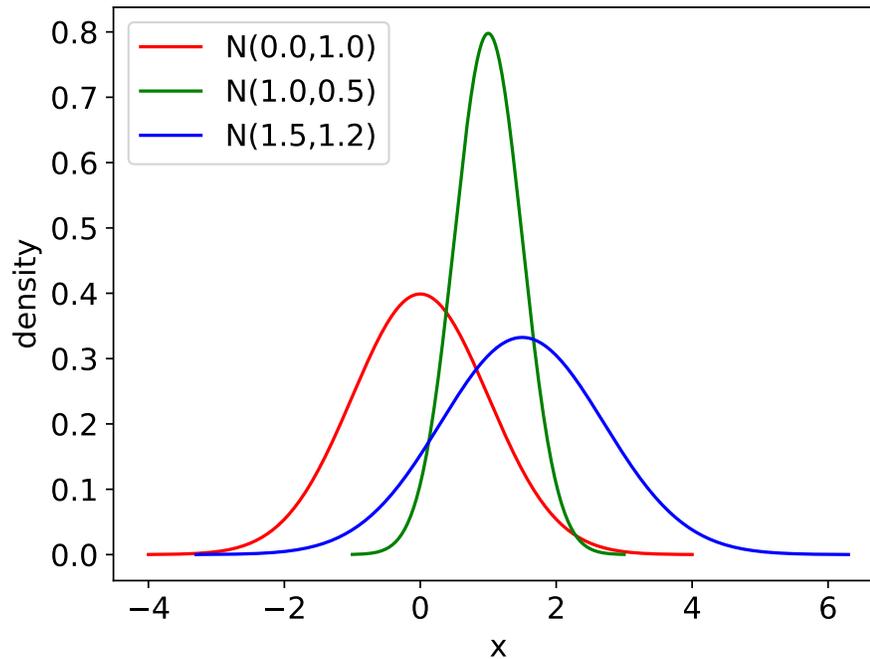
$$\phi(x) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right)$$



- More generally: a normal random variable with mean  $\mu$  and variance  $\sigma^2$  has density function

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

- Notation: “ $X \sim N(\mu, \sigma^2)$ ” means “ $X$  is a random variable with density function  $\phi_{\mu, \sigma^2}$ ”
- Fact: If  $X \sim N(0, 1)$  and  $Y = \mu + \sigma X$ , then  $Y \sim N(\mu, \sigma^2)$   
(Verify this using change-of-variable)



Q. What is the distribution function for  $X \sim \text{Unif}([0, 1])$ ?

## 1.9 Two important theorems

- Law of Large Numbers (LLN): If  $X_1, X_2, \dots$  is an infinite sequence of independent and identically distributed (i.i.d.) random variables with expectation  $\mu$ , then

$$\frac{1}{n} \sum_{i=1}^n X_i \longrightarrow \mu$$

as  $n \rightarrow \infty$

(We don't dwell upon the notions of convergence in this class)

- Central Limit Theorem (CLT): If  $X_1, X_2, \dots$  is an infinite sequence of independent and identically distributed (i.i.d.) random variables with expectation  $\mu$  and variance  $\sigma^2$ , then

$$\frac{\sum_{i=1}^n X_i - \mu}{\sigma\sqrt{n}} \longrightarrow N(0, 1)$$

as  $n \rightarrow \infty$

## 2 Review of linear algebra

### 2.1 Why linear algebra?

- Many machine learning methods represent data as vectors of numbers
- Many methods for statistical analysis is based on linear algebraic ideas (e.g., linearity)
- Descriptions and analyses of many machine learning methods use linear algebraic notations and concepts

### 2.2 Euclidean spaces

- Euclidean  $d$ -space, denoted  $\mathbb{R}^d$ , is the  $d$ -dimensional generalization of three-dimensional physical space
- A  $d$ -vector  $v \in \mathbb{R}^d$  is a  $d$ -tuple of real numbers

$$v = (v_1, \dots, v_d)$$

(We omit “ $d$ ” from “ $d$ -vector” when clear from context)

- The  $i$ -th component (a.k.a. entry) of  $v$  is  $v_i$
- Basic operations on  $d$ -vectors that produce  $d$ -vectors:

- Addition: for  $u, v \in \mathbb{R}^d$ ,

$$u + v = (u_1 + v_1, \dots, u_d + v_d) \in \mathbb{R}^d$$

- Scalar multiplication: for  $v \in \mathbb{R}^d$  and  $c \in \mathbb{R}$ ,

$$cv = (cv_1, \dots, cv_d) \in \mathbb{R}^d$$

- There is a special vector called the zero vector  $0 = (0, \dots, 0)$ 
  - Adding the zero vector to another vector  $v$  results in  $v$
  - Scaling the zero vector by a real number  $c$  results in the zero vector

- The norm (a.k.a. length) of a vector  $v \in \mathbb{R}^d$ , denoted by  $\|v\|$ , is

$$\|v\| = \sqrt{v_1^2 + \cdots + v_d^2}$$

– A unit vector is a vector with norm 1

- The inner product (a.k.a. dot product) between vectors  $u, v \in \mathbb{R}^d$ , denoted by  $u^\top v$  (or  $\langle u, v \rangle$ ), is

$$u^\top v = u_1 v_1 + \cdots + u_d v_d$$

– Interpretation:  $u^\top v = \|u\| \|v\| \cos(\theta)$  where  $\theta$  is the “angle” between  $u$  and  $v$

– Note:  $\|v\| = \sqrt{v^\top v}$

- Cauchy-Schwarz inequality: For any vectors  $u, v \in \mathbb{R}^d$ ,

$$u^\top v \leq \|u\| \|v\|,$$

with equality if and only if there is a real number  $c \in \mathbb{R}$  such that  $u = cv$

- Vectors  $u, v \in \mathbb{R}^d$  are orthogonal if  $u^\top v = 0$  (shorthand: “ $u \perp v$ ”)
  - A collection of vectors  $v^{(1)}, \dots, v^{(n)} \in \mathbb{R}^d$  is orthogonal if, for every  $i \neq j$ ,  $v^{(i)}$  and  $v^{(j)}$  are orthogonal
  - A collection of vectors is orthonormal if it is orthogonal and every vector in the collection is a unit vector
- Pythagorean theorem: If  $v^{(1)}, \dots, v^{(n)}$  is an orthogonal collection of vectors, then

$$\|v^{(1)} + \cdots + v^{(n)}\|^2 = \|v^{(1)}\|^2 + \cdots + \|v^{(n)}\|^2$$

Q. Show that the only vector with length zero is the zero vector.

Q. Show that the triangle inequality holds: for any  $u, v \in \mathbb{R}^d$ ,

$$\|u + v\| \leq \|u\| + \|v\|.$$

## 2.3 Linear dependence

- A linear combination of a finite collection of vectors  $v^{(1)}, \dots, v^{(n)} \in \mathbb{R}^d$  is an expression that multiplies each  $v^{(i)}$  by a real number  $c_i \in \mathbb{R}$ , and then adds up the results:

$$c_1v^{(1)} + \dots + c_nv^{(n)}$$

- A non-trivial linear combination of a finite collection of vectors  $v^{(1)}, \dots, v^{(n)} \in \mathbb{R}^d$  is a linear combination  $c_1v^{(1)} + \dots + c_nv^{(n)}$  where at least one of the  $c^{(i)}$  is non-zero
  - A collection of vectors is linearly dependent if there is a non-trivial linear combination of vectors from this collection that results in the zero vector
    - A collection of vectors that is not linearly dependent is said to be linearly independent
- Q. Suppose unit vectors  $v^{(1)}, \dots, v^{(n)}$  satisfy  $|\langle v^{(i)}, v^{(j)} \rangle| \leq 1/n$  for all  $i \neq j$ . Show that these vectors must be linearly independent.

## 2.4 Subspaces, dimension, and bases

- The span of a collection of vectors is the set of all linear combinations of any subset of vectors from this collection
- A subspace  $\mathcal{W}$  of  $\mathbb{R}^d$  is a collection of vectors from  $\mathbb{R}^d$  that is closed under addition and scalar multiplication and also contains the zero vector
  - $\mathbb{R}^d$  itself is a subspace of  $\mathbb{R}^d$
- The dimension of a subspace  $\mathcal{W}$ , written  $\dim(\mathcal{W})$ , is the largest number  $k$  such that  $\mathcal{W}$  contains a linearly independent set of  $k$  vectors
  - $\dim(\mathbb{R}^d) = d$
- A set of vector  $B$  from a subspace  $\mathcal{W}$  is a basis for  $\mathcal{W}$  if  $B$  is linearly independent and the span of  $B$  is  $\mathcal{W}$

- Every basis for a subspace  $\mathcal{W}$  has the same number of vectors, and that number is the dimension of the subspace
- It is often useful to order the vectors in a basis  $B = (b^{(1)}, \dots, b^{(k)})$ , and such an ordered set of vectors is called an ordered basis
- The standard (coordinate) basis for  $\mathbb{R}^d$  is the ordered basis  $(e^{(1)}, \dots, e^{(d)})$ , where  $e^{(i)}$  is the  $d$ -vector whose components are all zeros except for the  $i$ -th component, which has value one

## 2.5 Linear transformations and matrices

- A linear transformation  $T: \mathbb{R}^d \rightarrow \mathbb{R}^k$  between the Euclidean spaces  $\mathbb{R}^d$  and  $\mathbb{R}^k$  is a function that satisfies the following two properties:
  - Additivity:  $T(u + v) = T(u) + T(v)$  for any  $u, v \in \mathbb{R}^d$
  - Homogeneity:  $T(cv) = cT(v)$  for any  $v \in \mathbb{R}^d$  and  $c \in \mathbb{R}$

(Additivity & homogeneity = linearity)

- A  $k \times d$  matrix  $A$  is a tableaux of  $kd$  real numbers arranged in  $k$  rows and  $d$  columns

$$A = \begin{bmatrix} A_{1,1} & \cdots & A_{1,d} \\ \vdots & \ddots & \vdots \\ A_{k,1} & \cdots & A_{k,d} \end{bmatrix}$$

- The  $(i, j)$ -th component (a.k.a. entry) of  $A$  is  $A_{i,j}$
- We may regard  $A$  has an ordered collection of  $k$ -vectors  $a^{(1)}, \dots, a^{(d)}$ , one per column of  $A$ :

$$A = \begin{bmatrix} \uparrow & & \uparrow \\ a^{(1)} & \cdots & a^{(d)} \\ \downarrow & & \downarrow \end{bmatrix}$$

- The (matrix-vector) product of  $k \times d$  matrix  $A$  and  $d$ -vector  $x = (x_1, \dots, x_d)$ , written  $Ax$ , is the linear combination of columns  $a^{(1)}, \dots, a^{(d)} \in \mathbb{R}^k$  of  $A$  given by

$$x_1 a^{(1)} + \cdots + x_d a^{(d)}$$

- Caution: a matrix-vector product  $Ax$  only makes sense if the number of columns of  $A$  equals the number of components of  $x$

- The (matrix-matrix) product (a.k.a. matrix multiplication) of  $k \times d$  matrix  $A$  and  $d \times p$  matrix  $B$ , written  $AB$ , is the  $k \times p$  matrix whose  $i$ -th column is the matrix-vector product of  $A$  and the  $i$ -th column of  $B$

$$AB = [Ab^{(1)} \ \cdots \ Ab^{(p)}]$$

- Caution: a matrix-matrix product  $AB$  only makes sense if the number of columns of  $A$  equals the number of rows of  $B$
- Matrix-vector product can be viewed as a special case of matrix multiplication by pretending a  $d$ -vector is a  $d \times 1$  matrix
- Matrix multiplication is associative (i.e.,  $A(BC) = (AB)C$ ) and distributive (i.e.,  $A(B + C) = AB + AC$ ), but not commutative (i.e., it is possible that  $AB \neq BA$ )
- The transpose of  $k \times d$  matrix  $A$ , written  $A^T$  is the  $d \times k$  matrix whose  $(i, j)$ -th component is  $A_{j,i}$ 
  - What is the “meaning” of  $A^T$ ?
  - For every  $x \in \mathbb{R}^d$  and  $y \in \mathbb{R}^k$ , we have  $\langle Ax, y \rangle = \langle x, A^T y \rangle$
  - Relates certain angles in  $\mathbb{R}^d$  to certain angles in  $\mathbb{R}^k$
- Special matrix-matrix product: outer product of  $k$ -vector  $u \in \mathbb{R}^k$  and  $d$ -vector  $v \in \mathbb{R}^d$ , written  $uv^T$  (where  $v$  is treated as  $d \times 1$  matrix, so  $v^T$  is  $1 \times d$  matrix, a.k.a. row vector)

$$uv^T = \begin{bmatrix} u_1 \\ \vdots \\ u_k \end{bmatrix} [v_1 \ \cdots \ v_d] = \begin{bmatrix} u_1 v_1 & \cdots & u_1 v_d \\ \vdots & \ddots & \vdots \\ u_k v_1 & \cdots & u_k v_d \end{bmatrix}$$

(Result is a  $k \times d$  matrix)

- Q. Show that, for any  $k \times d$  matrix  $M$ , the transformation  $T: \mathbb{R}^d \rightarrow \mathbb{R}^k$  given by  $T(v) = Mv$  is a linear transformation.
- Q. Show that, for any linear transformation  $T: \mathbb{R}^d \rightarrow \mathbb{R}^k$ , there is a  $k \times d$  matrix  $M$  such that  $T(v) = Mv$  for all  $v \in \mathbb{R}^d$ .

## 2.6 Orthogonal complements and projectors

- For a  $k$ -dimensional subspace  $\mathcal{W}$  of  $\mathbb{R}^d$ , the orthogonal complement of  $\mathcal{W}$ , written  $\mathcal{W}^\perp$ , is the set of all vectors  $v$  that are orthogonal to every vector in  $\mathcal{W}$

$$\mathcal{W}^\perp = \{v \in \mathbb{R}^d : v \perp w \text{ for all } w \in \mathcal{W}\}$$

- Sometimes write “ $v \perp \mathcal{W}$ ” to mean “ $v \perp w$  for all  $w \in \mathcal{W}$ ”
- $\mathcal{W}^\perp$  is also a subspace of  $\mathbb{R}^d$
- For every  $k \times d$  matrix  $A$ :
  - Column space of  $A$ , denoted  $\text{CS}(A)$ , is the span of columns of  $A$
  - Nullspace of  $A$ , denoted  $\text{NS}(A)$ , is all  $x \in \mathbb{R}^d$  such that  $Ax = 0$
  - Row space of  $A$  is  $\text{CS}(A^\top)$ ; left nullspace of  $A$  is  $\text{NS}(A^\top)$
  - $\text{CS}(A^\top)$  and  $\text{NS}(A)$  are subspaces of  $\mathbb{R}^d$
  - $\text{CS}(A)$  and  $\text{NS}(A^\top)$  are subspaces of  $\mathbb{R}^k$
  - Rank of  $A$  is  $\dim(\text{CS}(A))$  and is also equal to  $\dim(\text{CS}(A^\top))$
  - $\text{rank}(A) + \dim(\text{NS}(A)) = d$ , and  $\text{rank}(A) + \dim(\text{NS}(A^\top)) = k$
  - $\text{CS}(A^\top)$  and  $\text{NS}(A)$  are orthogonal complements of each other
  - $\text{CS}(A)$  and  $\text{NS}(A^\top)$  are orthogonal complements of each other
- A projection operator (a.k.a. projector)  $P: \mathbb{R}^d \rightarrow \mathbb{R}^d$  is a linear transformation satisfying idempotency, i.e.,  $P(v) = P(P(v))$  for all  $v \in \mathbb{R}^d$
- For any subspace  $\mathcal{W}$  of  $\mathbb{R}^d$ , there is a projector  $P: \mathbb{R}^d \rightarrow \mathbb{R}^d$ , called the orthogonal projector (a.k.a. orthoprojector) to  $\mathcal{W}$ , such that

$$P(v) \in \mathcal{W} \quad \text{and} \quad v - P(v) \in \mathcal{W}^\perp$$

- If  $\mathcal{W} = \text{CS}(A)$  for a  $k \times d$  matrix  $A$ , then  $P(v) = AA^\dagger v$  for all  $v \in \mathbb{R}^d$ , where  $A^\dagger$  is Moore-Penrose pseudoinverse of  $A$ 
  - \* If  $\text{rank}(A) = d$ , then  $P(v) = A(A^\top A)^{-1}A^\top v$