

CS4618 Artificial Intelligence I

Today: Introduction to Probability (cont.)

Thomas Jansen

October 24th

Plans for Today

① Random Variables and Expectations

Introduction

Markov Inequality

② Conditional Probability

Introduction and the Law of Total Probability

③ Summary

Summary & Take Home Message

Introduction to Probability

Remember

- **probability space** is **sample space** (countable set) and **probability measure** (assigning probabilities to events)
- **consequence event and counter-event**
($\text{Prob}(A) + \text{Prob}(\bar{A}) = 1$)
- **consequence union bound** ($\text{Prob}(A_1 \cup A_2 \cup \dots \cup A_n) \leq \text{Prob}(A_1) + \text{Prob}(A_2) + \dots + \text{Prob}(A_n)$)
- **independence of events** (for all subsets
 $\text{Prob}(A_1 \cap A_2 \cap \dots \cap A_n) = \text{Prob}(A_1) \cdot \text{Prob}(A_2) \cdot \dots \cdot \text{Prob}(A_n)$)
- **random variable** maps elementary events to numbers
($X: \Omega \rightarrow \mathbb{R}$,
 $\forall x_i \in \mathbb{R}: \text{Prob}(X = x_i) = \text{Prob}(\{s \in \Omega: X(s) = x_i\})$)
- **expected value** gives average value, mean
($E(X) = \sum_{s \in \Omega} X(s) \cdot \text{Prob}(s)$)
- **Observation** Expected value may **not** correspond to possible

Direct Consequences

Theorem

Let X, Y be random variables on (Ω, Prob) , $a \in \mathbb{R}$.

- $E(X) = \sum_x x \cdot \text{Prob}(X = x)$
- $E(a \cdot X) = a \cdot E(X)$
- $E(X + Y) = E(X) + E(Y)$

Proof directly from definition □

Consider random variable $X: \Omega \rightarrow \{0, 1\}$

Observation $E(X) = \text{Prob}(X = 1)$

Proof $E(X) = 0 \cdot \text{Prob}(X = 0) + 1 \cdot \text{Prob}(X = 1)$
 $= \text{Prob}(X = 1)$ □

Direct Consequences (cont.)

Theorem

Let X, Y be independent random variables on (Ω, Prob) .

$$E(X \cdot Y) = E(X) \cdot E(Y)$$

Proof $E(X \cdot Y) = \sum_v v \cdot \text{Prob}(X \cdot Y = v)$

$$= \sum_x \sum_y (x \cdot y) \cdot \text{Prob}((X = x) \wedge (Y = y))$$

$$= \sum_x \sum_y (x \cdot y) \cdot \text{Prob}(X = x) \cdot \text{Prob}(Y = y)$$

$$= \left(\sum_x x \cdot \text{Prob}(X = x) \right) \cdot \left(\sum_y y \cdot \text{Prob}(Y = y) \right)$$

$$= E(X) \cdot E(Y)$$



Deviations from Expectation

Theorem (Markov Inequality)

$X \geq 0$ random variable, $s > 0$

$$\text{Prob}(X \geq s \cdot E(X)) \leq \frac{1}{s}$$

Proof Define $Y_t := \begin{cases} 1 & \text{if } X \geq t \\ 0 & \text{otherwise} \end{cases}$

Observe $X \geq t \cdot Y_t$

$$\begin{aligned} \Rightarrow E(X) &\geq E(t \cdot Y_t) = t \cdot E(Y_t) = t \cdot \text{Prob}(Y_t = 1) \\ &= t \cdot \text{Prob}(X \geq t) \end{aligned}$$

$$\Leftrightarrow \text{Prob}(X \geq t) \leq \frac{E(X)}{t}$$

$$\text{Set } t := s \cdot E(X) \Rightarrow \text{Prob}(X \geq s \cdot E(X)) \leq \frac{1}{s} \quad \square$$

Binomial Distribution

Consider n independent fair coin tosses ($\Omega = \{H, T\}^n$)

Consider $M :=$ number of coins landing heads

What is $\text{Prob}(M = i)$ and $E(M)$?

Observation $\text{Prob}(M = i) = \begin{cases} \binom{n}{i} \cdot 2^{-n} & \text{if } i \in \{0, 1, \dots, n\} \\ 0 & \text{otherwise} \end{cases}$

$$\Rightarrow E(M) = \sum_{i=0}^n i \cdot \binom{n}{i} \cdot 2^{-n} = ?$$

Define X_1, X_2, \dots, X_n with $X_i := \begin{cases} 1 & \text{if } i\text{-th coin lands heads} \\ 0 & \text{otherwise} \end{cases}$

Clearly $\forall i: E(X_i) = \frac{1}{2}$

Observation $M = \sum_{i=1}^n X_i$

$$\Rightarrow E(M) = E\left(\sum_{i=1}^n X_i\right) = \sum_{i=1}^n E(X_i) = \sum_{i=1}^n \text{Prob}(X_i = 1) = \frac{n}{2}$$

Geometric Distribution

Consider Sequence of independent and identical unfair coin tosses
for each coin $\text{Prob}(\text{coin lands heads}) = p$ with $p \in (0; 1)$

Consider

$T :=$ number of coin tosses until the first coin lands heads

What is $\text{Prob}(T = i)$ and $E(T)$?

Observation $\forall i \in \mathbb{N}_0: \text{Prob}(T = i) = p \cdot (1 - p)^{i-1}$

$$\begin{aligned} \Rightarrow E(T) &= \sum_{i=0}^{\infty} i \text{Prob}(T = i) = \sum_{i=0}^{\infty} ip(1-p)^{i-1} = \sum_{i=1}^{\infty} ip(1-p)^{i-1} \\ &= p \left(\sum_{i=1}^{\infty} \sum_{j=i}^{\infty} (1-p)^{j-1} \right) = p \left(\sum_{i=1}^{\infty} \left(\left(\sum_{j=1}^{\infty} (1-p)^{j-1} \right) - \left(\sum_{j=1}^{i-1} (1-p)^{j-1} \right) \right) \right) \\ &= p \left(\sum_{i=1}^{\infty} \left(\left(\sum_{j=0}^{\infty} (1-p)^j \right) - \left(\sum_{j=0}^{i-2} (1-p)^j \right) \right) \right) \\ &= p \cdot \left(\sum_{i=1}^{\infty} \frac{1}{1-(1-p)} - \frac{1-(1-p)^{i-1}}{1-(1-p)} \right) = p \cdot \left(\sum_{i=1}^{\infty} \frac{(1-p)^{i-1}}{p} \right) \\ &= \sum_{i=0}^{\infty} (1-p)^i = \frac{1}{1-(1-p)} = \frac{1}{p} \end{aligned}$$

Example Independent Fair Coin Tosses

Consider 100 independent fair coin tosses

Consider $M :=$ number of coins landing heads

We know $E(M) = \frac{1}{2} \cdot 100 = 50$

What is $\text{Prob}(M \geq 75)$?

Markov
$$\begin{aligned} \text{Prob}(M \geq 75) &= \text{Prob}\left(M \geq \frac{3}{2} \cdot 50\right) \\ &= \text{Prob}\left(M \geq \frac{3}{2} \cdot E(M)\right) \leq \frac{2}{3} \end{aligned}$$

Observation **weak bound**

$$\begin{aligned} \text{Prob}(M \geq 75) &= \sum_{i=75}^{100} \text{Prob}(M = i) = \sum_{i=75}^{100} \binom{100}{i} \cdot 2^{-100} \\ &= \frac{89,310,453,796,450,805,935,325}{316,912,650,057,057,350,374,175,801,344} < 0.000000282 \end{aligned}$$

Deviations from Expectation Revisted

Theorem (Chernoff Bounds)

Let $X_1, X_2, \dots, X_n: \Omega \rightarrow \{0, 1\}$ independent random variables with

$$\forall i \in \{1, 2, \dots, n\}: 0 < \text{Prob}(X_i = 1) < 1.$$

$$\text{Let } X := \sum_{i=1}^n X_i.$$

$$\forall \delta > 0: \text{Prob}(X > (1 + \delta) \cdot E(X)) < \left(\frac{e^\delta}{(1 + \delta)^{1 + \delta}} \right)^{E(X)}$$

$$\forall 0 < \delta < 1: \text{Prob}(X < (1 - \delta) \cdot E(X)) < e^{-E(X)\delta^2/2}$$

Proof Ideas

- Transform X to $e^{t \cdot X}$ for some $t > 0$.
- Apply Markov Inequality.
- Exploit independence of X_1, \dots, X_n .

Proof of Chernoff Bounds (1)

Let $t \in \mathbb{R}^+$.

Observe $\text{Prob}(X > (1 + \delta)\mathbb{E}(X)) = \text{Prob}\left(e^{tX} > e^{t(1+\delta)\mathbb{E}(X)}\right)$

Markov $\rightsquigarrow \text{Prob}(X > (1 + \delta)\mathbb{E}(X)) < \frac{\mathbb{E}(e^{tX})}{e^{t(1+\delta)\mathbb{E}(X)}}$

Remember $X = \sum_{i=1}^n X_i$

Thus $\mathbb{E}\left(e^{tX}\right) = \mathbb{E}\left(e^{t \cdot \sum_{i=1}^n X_i}\right) = \mathbb{E}\left(\prod_{i=1}^n e^{tX_i}\right)$

Remember X_i completely independent
 $\rightsquigarrow e^{tX_i}$ completely independent

Thus $\mathbb{E}\left(e^{tX}\right) = \prod_{i=1}^n \mathbb{E}\left(e^{tX_i}\right)$

Observe $\mathbb{E}\left(e^{tX_i}\right) = \text{Prob}(X_i = 0) \cdot 1 + \text{Prob}(X_i = 1) \cdot e^t$
 $= \text{Prob}(X_i = 1) \cdot e^t + 1 - \text{Prob}(X_i = 1)$
 $= 1 + \mathbb{E}(X_i) \cdot (e^t - 1) \leq e^{\mathbb{E}(X_i) \cdot (e^t - 1)}$

Proof of Chernoff Bounds (2)

We have $\text{Prob}(X > (1 + \delta)\mathbb{E}(X)) < \frac{\prod_{i=1}^n e^{\mathbb{E}(X_i) \cdot (e^t - 1)}}{e^{t(1+\delta)\mathbb{E}(X)}}$

$$= \frac{(e^t - 1) \cdot \sum_{i=1}^n \mathbb{E}(X_i)}{e^{t(1+\delta)\mathbb{E}(X)}} = \frac{e^{(e^t - 1) \cdot \mathbb{E}(X)}}{e^{t(1+\delta)\mathbb{E}(X)}}$$

We have $\delta > 0$, thus $t := \ln(1 + \delta)$ **feasible**

This yields $\text{Prob}(X > (1 + \delta)\mathbb{E}(X)) < \left(\frac{e^\delta}{(1+\delta)^{1+\delta}}\right)^{\mathbb{E}(X)}$ ✓

Analogously with $t := \ln\left(\frac{1}{1-\delta}\right)$

$$\text{Prob}(X < (1 - \delta)\mathbb{E}(X)) < \left(\frac{e^{-\delta}}{(1-\delta)^{1-\delta}}\right)^{\mathbb{E}(X)}$$

Observation $0 < \delta < 1 \Rightarrow (1 - \delta)^{1-\delta} > e^{-\delta + \delta^2/2}$

$$\Rightarrow \text{Prob}(X < (1 - \delta)\mathbb{E}(X)) < e^{-\mathbb{E}(X)\delta^2/2}$$



Example Independent Fair Coin Tosses Revisited

Consider 100 independent fair coin tosses
 $M :=$ number of coins landing heads ($E(M) = 50$)

What is $\text{Prob}(M \geq 75)$?

We know **Markov** $\text{Prob}(M \geq 75) \leq \frac{2}{3}$
truth $\text{Prob}(M \geq 75) < 0.000000282$

Can we apply Chernoff bounds?

Observe random variables X_1, X_2, \dots, X_{100} with

$$X_i := \begin{cases} 1 & i\text{-th coin lands heads} \\ 0 & \text{otherwise} \end{cases}$$

yields $M = \sum_{i=1}^n X_i$, X_1, X_2, \dots, X_n completely independent

$$\text{Prob}(M \geq 75) = \text{Prob}\left(M \geq \left(1 + \frac{1}{2}\right) \cdot E(M)\right) < \left(\frac{\sqrt{e}}{(3/2)^{(3/2)}}\right)^{50}$$

$$< 0.0044722$$

Taking Knowledge Into Account

Consider events A_1, A_2, \dots, A_n, B

What changes if we know that B occurred?

Clearly new probability distribution such that

$$\forall s \in \Omega \setminus B: \tilde{P}(s) = 0$$

Still need $\sum_{s \in \Omega \cap B} \tilde{P}(s) = 1$

Moreover relative probabilities inside B should not change

$$\forall A_i, A_j \subseteq B: \frac{\text{Prob}(A_i)}{\text{Prob}(A_j)} \stackrel{!}{=} \frac{\tilde{P}(A_i)}{\tilde{P}(A_j)}$$

Definition

For events A, B with $B \neq \emptyset$ define the **conditional probability of A given B** as $\text{Prob}(A | B) := \frac{\text{Prob}(A \cap B)}{\text{Prob}(B)}$.

Direct Consequences

Theorem (Law of Total Probability)

Let B_i with $i \in I$ be a *partition* of Ω .

$$\forall A \subseteq \Omega: \text{Prob}(A) = \sum_{i \in I} \text{Prob}(A | B_i) \cdot \text{Prob}(B_i)$$

Proof

$$\begin{aligned} \text{Prob}(A | B_i) \cdot \text{Prob}(B_i) &= \frac{\text{Prob}(A \cap B_i)}{\text{Prob}(B_i)} \cdot \text{Prob}(B_i) = \text{Prob}(A \cap B_i) \\ \sum_{i \in I} \text{Prob}(A | B_i) \cdot \text{Prob}(B_i) &= \sum_{i \in I} \text{Prob}(A \cap B_i) \\ &= \text{Prob}\left(\bigcup_{i \in I} (A \cap B_i)\right) = \text{Prob}\left(A \cap \bigcup_{i \in I} B_i\right) = \text{Prob}(A) \quad \square \end{aligned}$$

Corollary

$$\forall A, B \subseteq \Omega: \text{Prob}(A) \geq \text{Prob}(B) \cdot \text{Prob}(A | B)$$

Law of Total Probability for Random Variables

Clearly for random variable X , event B , $E(X | B)$ well defined

Theorem

Let B_i with $i \in I$ be a *partition* of Ω , X a random variable.

$$E(X) = \sum_{i \in I} E(X | B_i) \cdot \text{Prob}(B_i)$$

Proof Remember $E(X) = \sum_{s \in \Omega} X(s) \cdot \text{Prob}(s)$

$$\begin{aligned} E(X | B_i) \cdot \text{Prob}(B_i) &= \sum_{s \in \Omega} X(s) \cdot \text{Prob}(s | B_i) \cdot \text{Prob}(B_i) \\ &= \sum_{s \in \Omega} X(s) \cdot \text{Prob}(s \cap B_i) \end{aligned}$$

Observation $\forall s \in \Omega: \sum_{i \in I} \text{Prob}(s \cap B_i) = \text{Prob}(s)$

since B_i are a *partition*

$$\begin{aligned} \sum_{i \in I} E(X | B_i) \cdot \text{Prob}(B_i) &= \sum_{i \in I} \sum_{s \in \Omega} X(s) \cdot \text{Prob}(s \cap B_i) \\ &= \sum_{s \in \Omega} \left(X(s) \sum_{i \in I} \text{Prob}(s \cap B_i) \right) = \sum_{s \in \Omega} X(s) \text{Prob}(s) = E(X) \quad \square \end{aligned}$$

Summary & Take Home Message

Things to remember

- Markov inequality
- distributions: binomial, geometric
- Chernoff bounds
- conditional probability

Take Home Message

- Probability theory is not very complicated.
- Proving useful results is not very difficult.