

EECS126 Course Notes

Anmol Parande

Spring 2021 - Professor Thomas Courtade

Disclaimer: These notes reflect 126 when I took the course (Spring 2021). They may not accurately reflect current course content, so use at your own risk. If you find any typos, errors, etc, please raise an issue on the [GitHub repository](#).

Contents

1	Introduction to Probability	3
1.1	Conditional Probability	3
1.2	Independence	4
2	Random Variables and their Distributions	4
2.1	Distributions	5
2.2	Properties of Distributions	7
2.2.1	Expectation	7
2.2.2	Variance	9
2.3	Common Distributions	10
2.3.1	Discrete Distributions	10
2.3.2	Continuous Distributions	12
3	Concentration	13
3.1	Concentration Inequalities	14
3.2	Convergence	14
4	Information Theory	16

4.1	Quantifying Information	16
4.2	Source Coding	17
4.3	Channel Coding	19
5	Markov Chains	19
5.1	Properties of Markov Chains	20
5.1.1	Class Properties	21
5.2	Long-Term Behavior of Markov Chains	22

1 Introduction to Probability

Definition 1 A probability space is a triple (Ω, \mathcal{F}, P) where Ω is a set of objects called the sample space, \mathcal{F} is a family of subsets of Ω called events, and the probability measure $P : \mathcal{F} \rightarrow [0, 1]$.

One key assumption we make is that \mathcal{F} is a σ -algebra containing Ω , meaning that countably many complements, unions, and intersections of events in \mathcal{F} are also events in \mathcal{F} . The probability measure P must obey **Kolmogorov's Axioms**.

1. $\forall A \in \mathcal{F}, P(A) \geq 0$
2. $P(\Omega) = 1$
3. If $A_1, A_2, \dots \in \mathcal{F}$ and $\forall i \neq j, A_i \cap A_j = \emptyset$, then $P(\bigcup_{i \geq 1} A_i) = \sum_{i \geq 1} P(A_i)$

We choose Ω and \mathcal{F} to model problems in a way that makes our calculations easy.

Theorem 1

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

Theorem 2 (Inclusion-Exclusion Principle)

$$P\left(\bigcup_{i=1}^n A_i\right) = \sum_{k=1}^n (-1)^{k+1} \left(\sum_{1 \leq i_1 < \dots < i_k \leq n} P(A_{i_1} \cap \dots \cap A_{i_k}) \right)$$

Theorem 3 (Law of Total Probability) If A_1, A_2, \dots partition Ω (i.e. A_i are disjoint and $\bigcup A_i = \Omega$), then for event B ,

$$P(B) = \sum_i P(B \cap A_i)$$

1.1 Conditional Probability

Definition 2 If B is an event with $P(B) > 0$, then the conditional probability of A given B is

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

Intuitively, conditional probability is the probability of event A given that event B has occurred. In terms of probability spaces, it is as if we have taken (Ω, \mathcal{F}, P) and now have a probability measure $P(\cdot|C)$ belonging to the space $(\Omega, \mathcal{F}, P(\cdot|C))$.

Theorem 4 (Bayes Theorem)

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

1.2 Independence

Definition 3 Events A and B are independent if $P(A \cap B) = P(A)P(B)$

If $P(B) > 0$, then A, B are independent if and only if $P(A|B) = P(A)$. In other words, knowing B occurred gave no extra information about A .

Definition 4 If A, B, C with $P(C) > 0$ satisfy $P(A \cap B|C) = P(A|C)P(B|C)$, then A and B are conditionally independent given C .

Conditional independence is a special case of independence where A and B are not necessarily independent in the original probability space which has the measure P , but are independent in the new probability space conditioned on C with the measure $P(\cdot|C)$.

2 Random Variables and their Distributions

Definition 5 A random variable is a function $X : \Omega \rightarrow \mathbb{R}$ with the property $\forall \alpha \in \mathbb{R}, \{\omega \in \Omega : X(\omega) \leq \alpha\} \in \mathcal{F}$.

The condition in definition 5 is necessary to compute $P(X \leq \alpha), \forall \alpha \in \mathbb{R}$. This requirement also let us compute $P(X \in B)$ for most sets by leveraging the fact that \mathcal{F} is closed under complements, unions, and intersections. For example, we can also compute $P(X > \alpha)$ and $P(\alpha < X \leq \beta)$. In this sense, the property binds the probability space to the random variable.

definition 5 also implies that random variables satisfy particular algebraic properties. For example, if X, Y are random variables, then so are $X+Y, XY, X^p, \lim_{n \rightarrow \infty} X_n$, etc.

Definition 6 A discrete random variable is a random variable whose codomain is countable.

Definition 7 A continuous random variable is a random variable whose codomain is the real numbers.

Although random variables are defined based on a probability space, it is often most natural to model problems without explicitly specifying the probability space. This works so long as we specify the random variables and their distribution in a “consistent” way. This is formalized by the so-called [Kolmogorov Extension Theorem](#) but can largely be ignored.

2.1 Distributions

Roughly speaking, the distribution of a random variable gives an idea of the likelihood that a random variable takes a particular value or set of values.

Definition 8 The probability mass function (or distribution) of a random variable X is the frequency with which X takes on different values.

$$p_X : \mathcal{X} \rightarrow [0, 1] \text{ where } \mathcal{X} = \text{range}(X), \quad p_X(x) = \Pr \{X = x\}.$$

Note that $\sum_{x \in \mathcal{X}} p_X(x) = 1$ since $\bigcap_{x \in \mathcal{X}} \{w : X(w) = x\} = \Omega$.

Continuous random variables are largely similar to discrete random variables. One key difference is that instead of being described by a probability “mass”, they are instead described by a probability “density”.

Definition 9 The probability density function (distribution) of a continuous random variable describes the density by which a random variable takes a particular value.

$$f_X : \mathbb{R} \rightarrow [0, \infty) \text{ where } \int_{-\infty}^{\infty} f_X(x) dx = 1 \text{ and } \Pr \{X \in B\} = \int_B f_X(x) dx$$

Observe that if a random variable X is continuous, then the probability that it takes on a particular value is zero.

$$\Pr \{X = x\} = \lim_{\delta \rightarrow 0} \Pr \{x \leq X \leq x + \delta\} = \lim_{\delta \rightarrow 0} \int_x^{x+\delta} f_X(u) du = \int_x^x f_X(u) du = 0$$

Definition 10 The cumulative distribution function (CDF) gives us the probability of a random variable X being less than or equal to a particular value.

$$F_X : \mathbb{R} \rightarrow [0, 1], \quad F_X(x) = \Pr \{X \leq x\}$$

Note that by the Kolomogorov axioms, F_X must satisfy three properties:

1. F_X is non-decreasing.
2. $\lim_{x \rightarrow 0} F_X(x) = 0$ and $\lim_{x \rightarrow \infty} F_X(x) = 1$.
3. F_X is right continuous.

It turns out that if we have any function F_X that satisfies these three properties, then it is the CDF of some random variable on some probability space. Note that $F_X(x)$ gives us an alternative way to define continuous random variables. If $F_X(x)$ is absolutely continuous, then it can be expressed as

$$F_X(x) = \int_{-\infty}^x f_X(x) dx$$

for some non-negative function $f_X(x)$, and this is the PDF of a continuous random variable.

Often, when modeling problems, there are multiple random variables that we want to keep track of.

Definition 11 *If X and Y are random variables on a common probability space (Ω, \mathcal{F}, P) , then the joint distribution (denoted $p_{XY}(x, y)$ or $f_{XY}(x, y)$) describes the frequencies of joint outcomes.*

Note that it is possible for X to be continuous and Y to be discrete (or vice versa).

Definition 12 *The marginal distribution of a joint distribution is the distribution of a single random variable.*

$$p_X(x) = \sum_y p_{XY}(x, Y = y), \quad f_X(x) = \int_{-\infty}^{\infty} f_{XY}(x, y) dy$$

Definition 13 *Two random variables X and Y are independent if their joint distribution is the product of the marginal distributions.*

Just like independence, we can extend the notion of conditional probability to random variables.

Definition 14 *The conditional distribution of X given Y captures the frequencies of X given we know the value of Y .*

$$p_{X|Y}(x|y) = \frac{P_{XY}(x, y)}{p_Y(y)}, \quad f_{X|Y}(x|y) = \frac{f_{XY}(x, y)}{f_Y(y)}$$

Often, we need to combine or transform several random variables. A derived distribution is the obtained by arithmetic of several random variables or applying a function to several (or many) random variables. Since the CDF of a distribution essentially defines that random variable, it can often be easiest to work backwards from the CDF to the PDF or PMF. In the special case where we want to find $Y = g(X)$ for a function g .

$$F_y(y) = \Pr\{Y \leq y\} = \Pr\{g(x) \leq y\} = \Pr\{X \in g^{-1}([-\infty, y])\}, \quad g^{-1}(y) = \{x : g(x) = y\}.$$

Another special case of a derived distribution is when adding random variables together.

Theorem 5 *The resulting distribution of a sum of two independent random variables is the convolution of the distributions of the two random variables.*

$$p_{X+Y}(z) = \sum_{k=-\infty}^{\infty} p_X(k)p_Y(z-k), \quad f_{X+Y}(z) = \int_{-\infty}^{\infty} f_X(x)f_Y(z-x)dx$$

2.2 Properties of Distributions

2.2.1 Expectation

Definition 15 *The expectation of a random variable describes the center of a distribution,*

$$\mathbb{E}[X] = \sum_{x \in \mathcal{X}} xp_X(x), \quad \mathbb{E}[X] = \int_{-\infty}^{\infty} xf_X(x)dx$$

provided the sum or integral converges.

Expectation has several useful properties. If we want to compute the expectation of a function of a random variable, then we can use the law of the unconscious statistician.

Theorem 6 (Law of the Unconscious Statistician)

$$\mathbb{E}[g(X)] = \sum_{x \in \mathcal{X}} g(x)p_X(x), \quad \mathbb{E}[g(X)] = \int_{-\infty}^{\infty} g(x)f_X(x)dx$$

Another useful property is its linearity.

$$\mathbb{E}[aX + bY] = a\mathbb{E}[X] + b\mathbb{E}[Y], \quad \forall a, b \in \mathbb{R}.$$

Sometimes it can be difficult to compute expectations directly. For discrete distributions, we can use the tail-sum formula.

Theorem 7 (Tail Sum) For a non-negative integer random variable,

$$\mathbb{E}[X] = \sum_{k=1}^{\infty} \Pr\{X \geq k\}.$$

When two random variables are independent, expectation has some additional properties.

Theorem 8 If X and Y are independent, then

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

We can apply expectations to conditional distributions as well.

Definition 16 The conditional expectation of a conditional distribution is given by

$$\mathbb{E}[X|Y=y] = \sum_{x \in \mathcal{X}} x p_{X|Y}(x|y), \quad \mathbb{E}[X|Y=y] = \int_{-\infty}^{\infty} x f_{X|Y}(x, y) dx$$

Notice that $\mathbb{E}[X|Y]$ is a function of the random variable Y , meaning we can apply theorem 6.

Theorem 9 (Tower Property) For all functions f ,

$$\mathbb{E}[f(Y)X] = \mathbb{E}[f(Y)\mathbb{E}[X|Y]]$$

If we apply theorem 9 to the function $f(Y) = 1$, then we can see that $\mathbb{E}[\mathbb{E}[X|Y]] = \mathbb{E}[X]$.

Earlier, we saw that we find a derived distribution by transforming and combining random variables. Sometimes, we don't need to actually compute the distribution, but only some of its properties.

Definition 17 The n th moment of a random variable is $\mathbb{E}[X^n]$.

It turns out that we can encode the moments of a distribution into the coefficients of a special power series.

Definition 18 *The moment generating function of a random variable X is given by $M_X(t) = \mathbb{E}[e^{tX}]$.*

Notice that if we apply the power series expansion of e^{tX} , we see that

$$M_X(t) = \sum_{n=0}^{\infty} \frac{t^n}{n!} \mathbb{E}[X^n].$$

Thus the n th moment is encoded in the coefficients of the power series and we can retrieve them by taking a derivative:

$$\mathbb{E}[X^n] = \frac{d^n}{dt^n} M_X(t).$$

Another interesting point to notice is that for a continuous random variable

$$M_X(t) = \int_{-\infty}^{\infty} f_X(x) e^{tx} dx$$

is the Laplace transform of the distribution over the real line, and for a discrete random variable,

$$M_X(t) = \sum_{x=-\infty}^{\infty} p_X(x) e^{tx}$$

is the Z-transform of the distribution evaluated along the curve at e^{-t} .

Theorem 10 *If the MGF of a function exists, then it uniquely determines the distribution.*

This provides another way to compute the distribution for a sum of random variables because we can just multiply their MGF.

2.2.2 Variance

Definition 19 *The variance of a discrete random variable X describes its spread around the expectation and is given by*

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

Theorem 11 *When two random variables X and Y are independent, then*

$$\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y).$$

Definition 20 The covariance of two random variables describes how much they depend on each other and is given by

$$\text{Cov}(X, Y) = \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])] = \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y].$$

If $\text{Cov}(X, Y) = 0$ then X and Y are uncorrelated.

Definition 21 The correlation coefficient gives a single number which describes how random variables are correlated.

$$\rho(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)}\sqrt{\text{Var}(Y)}}.$$

Note that $-1 \leq \rho \leq 1$. Just as expectation can change when we know additional information, so can variance.

Definition 22 Conditional Variance is the variance of X given the value of Y .

$$\text{Var}(X|Y = y) = \mathbb{E}[(X - \mathbb{E}[X|Y = y])^2|Y = y] = \mathbb{E}[X^2|Y = y] - \mathbb{E}[X|Y = y]^2$$

Conditional variance is a random variable just as expectation is.

Theorem 12 (Law of Total Variance)

$$\text{Var}(X) = \mathbb{E}[\text{Var}(X|Y)] + \text{Var}(\mathbb{E}[X|Y])$$

The second term in the law of total variance ($\text{Var}(\mathbb{E}[X|Y])$) can be interpreted as on average, how much uncertainty there is in X given we know Y .

2.3 Common Distributions

2.3.1 Discrete Distributions

Definition 23 X is uniformly distributed when each value of X has equal probability.

$$X \sim \text{Uniform}(\{1, 2, \dots, n\}) \implies p_X(x) = \begin{cases} \frac{1}{n} & x = 1, 2, \dots, n, \\ 0 & \text{else.} \end{cases}$$

Definition 24 X is a Bernoulli random variable if it is either 0 or 1 with $p_X(1) = p$.

$$X \sim \text{Bernoulli}(p) \implies p_X(x) = \begin{cases} 1-p & x=0, \\ p & x=1, \\ 0 & \text{else.} \end{cases}$$

$$\mathbb{E}[X] = p \quad \text{Var}(X) = (1-p)p$$

Bernoulli random variables are good for modeling things like a coin flip where there is a probability of success. Bernoulli random variables are frequently used as indicator random variables $\mathbb{1}_A$ where

$$\mathbb{1}_A = \begin{cases} 1 & \text{if } A \text{ occurs,} \\ 0 & \text{else.} \end{cases}$$

When paired with the linearity of expectation, this can be a powerful method of computing the expectation of something.

Definition 25 X is a Binomial random variable when

$$X \sim \text{Binomial}(n, p) \implies p_X(x) = \begin{cases} \binom{n}{x} p^x (1-p)^{n-x} & x=0, 1, \dots, n \\ 0 & \text{else.} \end{cases}$$

$$\mathbb{E}[X] = np \quad \text{Var}(X) = np(1-p)$$

A binomial random variable can be thought of as the number of successes in n trials. In other words,

$$X \sim \text{Binomial}(n, p) \implies X = \sum_{i=1}^n X_i, \quad X_i \sim \text{Bernoulli}(p).$$

By construction, if $X \sim \text{Binomial}(n, p)$ and $Y \sim \text{Binomial}(m, p)$ are independent, then $X + Y \sim \text{Binomial}(m+n, p)$.

Definition 26 A Geometric random variable is distributed as

$$X \sim \text{Geom}(p) \implies p_X(x) = \begin{cases} p(1-p)^{x-1} & x=1, 2, \dots \\ 0 & \text{else.} \end{cases}$$

$$\mathbb{E}[X] = \frac{1}{p} \quad \text{Var}(X) = \frac{1-p}{p^2}$$

Geometric random variables are useful for modeling the number of trials required before the first success. In other words,

$$X \sim \text{Geom}(p) \implies X = \min\{k \geq 1 : X_k = 1\} \text{ where } X_i \sim \text{Bernoulli}(p).$$

A useful property of geometric random variables is that they are memoryless:

$$\Pr\{X = K + M | X > k\} = \Pr\{X = M\}.$$

Definition 27 A Poisson random variable is distributed as

$$X \sim \text{Poisson}(\lambda) \implies p_X(x) = \begin{cases} \frac{\lambda^x e^{-\lambda}}{x!} & x = 0, 1, \dots \\ 0 & \text{else.} \end{cases}$$

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

Poisson random variables are good for modeling the number of arrivals in a given interval. Suppose you take a given time interval and divide it into n chunks where the probability of arrival in chunk i is $X_i \sim \text{Bernoulli}(p_n)$. Then the total number of arrivals $X_n = \sum_{i=1}^n X_i$ is distributed as a Binomial random variable with expectation $np_n = \lambda$. As we increase n to infinity but keep λ fixed, we arrive at the poisson distribution.

A useful fact about Poisson random variables is that if $X \sim \text{Poisson}(\lambda)$ and $Y \sim \text{Poisson}(\mu)$ are independent, then $X + Y \sim \text{Poisson}(\lambda + \mu)$.

2.3.2 Continuous Distributions

Definition 28 A continuous random variable is uniformly distributed when the pdf of X is constant over a range.

$$X \sim \text{Uniform}(a, b) \implies f_X(x) = \begin{cases} \frac{1}{b-a} & a \leq x \leq b, \\ 0 & \text{else.} \end{cases}$$

The CDF of a uniform distribution is given by

$$F_X(x) = \begin{cases} 0, & x < a, \\ \frac{x-a}{b-a}, & x \in [a, b) \\ 1, & x \geq b. \end{cases}$$

Definition 29 A continuous random variable is exponentially distributed when its pdf is given by

$$X \sim \text{Exp}(\lambda) \implies f_X(x) = \begin{cases} \lambda e^{-\lambda x} & x \geq 0, \\ 0 & \text{else.} \end{cases}$$

Exponential random variables are the only continuous random variable to have the memoryless property:

$$\Pr\{X > t + s | X > s\} = \Pr\{X > t\}, \quad t \geq 0.$$

The CDF of the exponential distribution is given by

$$F_X(x) = \lambda \int_0^x e^{-\lambda u} du = 1 - e^{-\lambda x}$$

Definition 30 X is a Gaussian Random Variable with mean μ and variance σ^2 (denoted $X \sim \mathcal{N}(\mu, \sigma^2)$) if it has the PDF

$$f_X(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

The standard normal is $X \sim \mathcal{N}(0, 1)$, and it has the CDF

$$\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{u^2}{2}} du$$

There is no closed form for $\Phi(x)$. It turns out that every normal random variable can be transformed into the standard normal (i.e. $\frac{X-\mu}{\sigma} \sim \mathcal{N}(0, 1)$). Some facts about Gaussian random variables are

1. If $X \sim \mathcal{N}(\mu_x, \sigma_x^2)$, $Y \sim \mathcal{N}(\mu_y, \sigma_y^2)$ are independent, then $X + Y \sim \mathcal{N}(\mu_x + \mu_y, \sigma_x^2 + \sigma_y^2)$.
2. If X, Y are independent and $(X + Y), (X - Y)$ are independent, then both X and Y are Gaussian with the same variance.

3 Concentration

In real life, for the most part, we can't compute probabilities in closed form. Instead, we either bound them, or we want to show that $P(A) \approx 0$ or $P(A) \approx 1$.

3.1 Concentration Inequalities

Theorem 13 (Markov's Inequality) *For a non-negative random variable X ,*

$$\Pr \{X \geq t\} \leq \frac{\mathbb{E}[X]}{t}, \quad t \geq 0.$$

Theorem 14 (Chebyshev's Inequality) *If X is a random variable, then*

$$\Pr \{|X - \mathbb{E}[X]| \geq t\} \leq \frac{\text{Var}(X)}{t^2}.$$

Intuitively, theorem 14 gives a “better” bound than theorem 13 because it incorporates the variance of the random variable. Using this idea, we can define an even better bound that incorporates information from all moments of the random variable.

Definition 31 (Chernoff Bound) *For a random variable X and $a \in \mathbb{R}$,*

$$\Pr \{X \geq a\} \leq \frac{\mathbb{E}[e^{tX}]}{e^{ta}} = e^{-ta} M_x(t).$$

After computing the Chernoff bound for a general t , we can then optimize over it to compute the best bound possible.

3.2 Convergence

The idea of convergence brings the mathematical language of limits into probability. The fundamental question we want to answer is given random variables X_1, X_2, \dots , what does it mean to compute

$$\lim_{n \rightarrow \infty} X_n.$$

This question is not as straightforward as it seems because random variables are functions, and there are many ways to define the convergence of functions.

Definition 32 *A sequence of random variables converges almost surely to X if*

$$P \left(\lim_{n \rightarrow \infty} X_n = X \right) = 1$$

One result of almost sure convergence deals with deviations around the mean of many samples.

Theorem 15 (Strong Law of Large Numbers) *If X_1, X_2, \dots, X_n are independently and identically distributed to X where $\mathbb{E}[X] < \infty$, then $\frac{1}{n} \sum_i X_i$ converges almost surely to $\mathbb{E}[X]$.*

The strong law tells us that for any observed realization, there is a point after which there are no deviations from the mean.

Definition 33 *A sequence of random variables converges in probability if*

$$\forall \epsilon > 0, \quad \lim_{n \rightarrow \infty} P(|X_n - X| > \epsilon) = 0$$

Convergence in probability can help us formalize the intuition that we have which says probability is the frequency with which an even happens over many trials of an event.

Theorem 16 (Weak Law of Large Numbers) *Let X_1, X_2, \dots, X_n be independently and identically distributed according to X , and let $M_n = \frac{1}{n} \sum X_i$. Then for $\epsilon > 0$,*

$$\lim_{n \rightarrow \infty} \Pr \{ |M_n - \mathbb{E}[X]| > \epsilon \} = 0.$$

It tells us that the probability of a deviation of ϵ from the true mean will go to 0 in the limit, but we can still observe these deviations. Nevertheless, the weak law helps us formalize our intuition about probability. If X_1, X_2, \dots, X_n are independently and identically distributed according to X , then we can define the empirical frequency

$$F_n = \frac{\sum \mathbb{1}_{X_i \in B}}{n} \implies \mathbb{E}[F_n] = P(X \in B).$$

By theorem 16,

$$\lim_{n \rightarrow \infty} \Pr \{ |F_n - P(X \in B)| > \epsilon \} = 0,$$

meaning over many trials, the empirical frequency is equal to the probability of the event, matching intuition.

Definition 34 *A sequence of random variables converges in distribution if*

$$\lim_{n \rightarrow \infty} F_{X_n}(x) = F_x(x).$$

An example of convergence in distribution is the central limit theorem.

Theorem 17 (Central Limit Theorem) *If X_1, X_2, \dots are independently and identically distributed according to X with $\text{Var}(X) = \sigma^2$ and $\mathbb{E}[X] = \mu$, then*

$$\lim_{n \rightarrow \infty} P\left(\frac{\sum_{i=1}^n X_i - n\mu}{\sigma\sqrt{n}} \leq x\right) = \Phi(x)$$

In other words, a sequence of random variables converges in distribution to a normal distribution with variance σ^2 and mean μ .

These notions of convergence are not identical, and they do not necessarily imply each other. It is true that almost sure convergence implies convergence in probability, and convergence in probability implies convergence in distribution, but the implication is only one way.

4 Information Theory

Information Theory is a field which addresses two questions

1. **Source Coding:** How many bits do I need to losslessly represent an observation.
2. **Channel Coding:** How reliably and quickly can I communicate a message over a noisy channel.

4.1 Quantifying Information

Intuitively, for a PMF of a discrete random variable, the surprise associated with a particular realization is $-\log p_X(x)$ since less probable realizations are more surprising. With this intuition, we can try and quantify the “expected surprise” of a distribution.

Definition 35 *For a Discrete Random Variable $X \sim p_X$, the Entropy of X is given by*

$$H(x) = \mathbb{E}[-\log_2 p_X(x)] = -\sum_{x \in \mathcal{X}} p_X(x) \log_2 p_X(x).$$

Alternative interpretations of entropy are the average uncertainty and how random X is. Just like probabilities, we can define both joint and conditional entropies.

Definition 36 For Discrete Random Variables X and Y , the joint entropy is given by

$$H(X, Y) = \mathbb{E}[-\log_2 p_{XY}(x, y)] = - \sum_{x, y \in \mathcal{X} \times \mathcal{Y}} p_{XY}(x, y) \log_2 p_{XY}(x, y).$$

Definition 37 For Discrete Random Variable X and Y , the conditional entropy is given by

$$H(Y|X) = \mathbb{E}[-\log_2 p_{Y|X}(y|x)] = - \sum_{x \in \mathcal{X}} p_X(x) H(Y|X = x).$$

Conditional entropy has a natural interpretation which is that it tells us how surprised we are to see $Y = y$ given that we know $X = x$. If X and Y are independent, then $H(Y) = H(Y|X)$ because realizing X gives no additional information about Y .

Theorem 18 (Chain Rule of Entropy)

$$H(X, Y) = H(X) + H(Y|X).$$

In addition to knowing how much our surprise changes for a random variable when we observe a different random variable, we can also quantify how much additional information observing a random variable gives us about another.

Definition 38 For random variables X and Y , the mutual information is given by

$$I(X; Y) = H(X) - H(X|Y) = H(Y) - H(Y|X).$$

4.2 Source Coding

Source coding deals with finding the minimal number of bits required to represent data. This is essentially the idea of lossless compression. In this case, our message is the sequence of realizations of independently and identically distributed random variables $(X_i)_{i=1}^n \sim p_X$. The probability of observing a particular sequence is then

$$P(x_1, x_2, \dots, x_n) = \prod_{i=1}^n p_X(x_i).$$

Theorem 19 (Asymptotic Equipartition Property) *If we have a sequence of independently and identically distributed random variables $(X_i)_{i=1}^n \sim p_X$, then $-\frac{1}{n} \log P(x_1, x_2, \dots, x_n)$ converges to $H(X)$ in probability.*

theorem 19 tells us that with overwhelming probability, we will observe a sequence that is assigned probability $2^{-nH(X)}$. Using this idea, we can define a subset of possible observed sequences that in the limit, our observed sequence must belong to with overwhelming probability.

Definition 39 *For a fixed $\epsilon > 0$, for each $n \geq 1$, the typical set is given by*

$$A_\epsilon^{(n)} = \{(x_1, x_2, \dots, x_n) : P(x_1, x_2, \dots, x_n) \geq 2^{-n(H(X)+\epsilon)}\}.$$

Two important properties of the typical set are that

1. $\lim_{n \rightarrow \infty} P((x_1, x_2, \dots, x_n) \in A_\epsilon^{(n)}) = 1$
2. $|A_\epsilon^{(n)}| \geq 2^{n(H(X)+\epsilon)}$

The typical set gives us an easy way to do source coding. If I have N total objects, then I only need $\log N$ bits to represent each object, so I can define a simple protocol which is

1. If $(x_i)_{i=1}^n \in A_{\frac{\epsilon}{2}}^{(n)}$, then describe them using the $\log |A_{\frac{\epsilon}{2}}^{(n)}| \leq n(H(X) + \frac{\epsilon}{2})$ bits
2. If $(x_i)_{i=1}^n \notin A_{\frac{\epsilon}{2}}^{(n)}$, then describe them naively with $n \log |\mathcal{X}|$ bits.

This makes the average number of bits required to describe a message

$$\begin{aligned} \mathbb{E}[\text{\# of Bits}] &\leq n(H(X) + \frac{\epsilon}{2})P((x_i)_{i=1}^n \in A_{\frac{\epsilon}{2}}^{(n)}) + n \log |\mathcal{X}|P((x_i)_{i=1}^n \notin A_{\frac{\epsilon}{2}}^{(n)}) \\ &\leq n(H(X) + \frac{\epsilon}{2}) + n \frac{\epsilon}{2} \leq n(H(X) + \epsilon) \end{aligned}$$

This is the first half of a central result of source coding.

Theorem 20 (Source Coding Theorem) *If $(X_i)_{i=1}^n \sim p_X$ are a sequence of independently and identically distributed random variable, then for any $\epsilon > 0$ and n sufficiently large, we can represent $(X_i)_{i=1}^n$ using fewer than $n(H(X) + \epsilon)$ bits. Conversely, we can not losslessly represent $(X_i)_{i=1}^n$ using fewer than $nH(X)$ bits.*

This lends a new interpretation of the entropy $H(X)$: it is the average number of bits required to represent X .

4.3 Channel Coding

Whereas source coding deals with encoding information, channel coding deals with transmitting it over a noisy channel. In general, we have a message M , and encoder, a channel, and a decoder as in fig. 1.



Figure 1: Channel Coding

Each channel can be described by a conditional probability distribution $p_{Y|X}(y|x)$ for each time the channel is used.

Definition 40 For a channel described by $p_{Y|X}$, the capacity is given by

$$C = \max_{p_X} I(X; Y).$$

In words, the capacity describes the maximum mutual information between the channel input and output.

Definition 41 Suppose we use the channel n times to send a message that takes on average $H(m)$ bits to encode, then the rate of the channel is

$$R = \frac{H(M)}{n}$$

Theorem 21 Channel Coding Theorem For a channel described by $p_{Y|X}$ and $\epsilon > 0$ and $R < C$, for all n sufficiently large, there exists a rate R communication scheme that achieves a probability of error less than ϵ . If $R > C$, then the probability of error converges to 1 for any communication scheme.

5 Markov Chains

Definition 42 A random/stochastic process is a sequence of random variables $(X_n)_{n \geq 0}$.

The random variables in a stochastic process do not have to be independently and identically distributed. In fact, if they are not, then we can get additional modeling power.

Definition 43 $(X_n)_{n \geq 0}$ is a Markov Chain if each random variable X_i takes values in a discrete set S (the state space), and,

$$\forall n \geq 0, i, j \in S, \Pr \{X_{n+1} = j | X_n = i, \dots, X_0 = x_0\} = \Pr \{X_{n+1} = j | X_n = i\}$$

In words, a Markov Chain is a sequence of random variables satisfying the Markov Property where probability of being in a state during the next time step only depends on the current state.

Definition 44 A temporally homogenous Markov Chain is one where the transition probabilities $\Pr \{X_{n+1} = j | X_n = i\} = p_{ij}$ for all $i, j \in S$ and $n \geq 0$.

Temporally Homogenous Markov Chains don't change their transition probabilities over time. Since the p_{ij} are conditional probabilities, they must satisfy

1. $\forall i, j \in S, p_{ij} \geq 0$
2. $\forall i \in S, \sum_{j \in S} p_{ij} = 1$

Definition 45 The transition matrix of a Markov Chain is a matrix P where the ij th entry $P_{ij} = p_{ij}$ for all $i, j \in S$.

The transition matrix encodes the one-step transition probabilities of the Markov Chain.

Theorem 22 (Chapman-Kolmogorov Equation) The n -step transition probabilities (i.e starting in i and ending in j n steps later) of the Markov Chain are given by $p_{ij}^{(n)} = P_{ij}^n$.

5.1 Properties of Markov Chains

Definition 46 If $\exists n \geq 1$ such that $p_{ij}^{(n)} \neq 0$, then j is accessible from i , and we write $i \rightarrow j$.

Definition 47 States i and j communicate with each other when $i \rightarrow j$ and $j \rightarrow i$. We write this as $i \leftrightarrow j$.

By convention, we say that $i \leftrightarrow i$. It turns out that \leftrightarrow is an equivalence relation on the state space S . An equivalence relation means that

1. $\forall i \in S, i \leftrightarrow i$

2. $\forall i, j \in S, i \leftrightarrow j \Leftrightarrow j \leftrightarrow i$
3. $\forall i, j, k \in S, i \leftrightarrow k, k \leftrightarrow j \Rightarrow i \leftrightarrow j$

This means that \leftrightarrow partitions the state-space S into equivalence classes (i.e classes of communicating states).

Definition 48 *A Markov Chain is irreducible if S is the only class.*

5.1.1 Class Properties

A class property is a property where if one element of a class has the property, all elements of the class have the property. Markov Chains have several of these properties which allow us to classify states.

Definition 49 *A state $i \in S$ is recurrent if given that $X_0 = i$, the process revisits state i with probability 1.*

Definition 50 *A state $i \in S$ is transient if it is not recurrent.*

Recurrence means that we will visit a state infinitely often in the future if we start in that state, while transience means we will only visit the state finitely many times. We can further break recurrence down if we define $T_i = \min_{n \geq 1} n$ such that $X_n = i$ (the first time the chain enters state i).

Definition 51 *State i is positive recurrent if it is recurrent and $\mathbb{E}[T_i | X_0 = i]$ is finite.*

Definition 52 *State i is null recurrent if it is recurrent and $\mathbb{E}[T_i | X_0 = i]$ is infinite.*

Positive recurrence means we visit a recurrent state so frequently that we spend a positive fraction of time in that state. Null recurrence means we visit a recurrent state so infrequently (but still infinitely many times) that we spend virtually no time in that state.

Theorem 23 *Every irreducible finite state Markov Chain is positive recurrent.*

Definition 53 *For a state $i \in S$, we define the period of the state to be*

$$\text{period}(i) = \text{GCD}\{n \geq 1 : p_{ii}^{(n)} > 0\}.$$

If we start in state i , then revisits to i only occur at integer multiples of the period.

Definition 54 *An irreducible markov chain is aperiodic if any state has period 1.*

All of the above properties are class properties.

5.2 Long-Term Behavior of Markov Chains

Since the p_{ij} completely characterize the Markov Chain, we can also describe what happens to the chain in the limit.

Definition 55 A probability distribution π over the states is a stationary distribution if $\pi = \pi P$

It is called a stationary distribution because the distribution over states is invariant with time.

Theorem 24 (Big Theorem for Markov Chains) Let $(X_n)_{n \geq 0}$ be an irreducible Markov Chain. Then one of the following is true.

1. Either all states are transient, or all states are null recurrent, and no stationary distribution exists, and $\lim_{n \rightarrow \infty} p_{ij}^{(n)}$.
2. All states are positive recurrent and the stationary distribution exists, is unique, and satisfies

$$\pi_j = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=0}^n P_{ij}^{(k)} = \frac{1}{\mathbb{E}[T_j | X_0 = j]}.$$

If the Markov Chain is periodic, then $\lim_{n \rightarrow \infty} p_{ij}^{(n)} = \pi_j$