

Problem Set 5

published: 19 May 2022
discussion: 27 May 2022

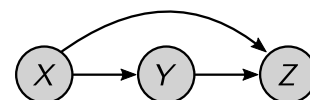
Data Compression With And Without Deep Probabilistic Models

Prof. Robert Bamler, University of Tuebingen

Course materials available at <https://robamler.github.io/teaching/compress22/>

Problem 5.1: Conditional Independence

In last week's lecture, we learned that every probability distribution P satisfies the so-called chain rule of probability theory. For example, for any three random variables X , Y , and Z , we can always factorize their joint probability distribution as follows (see also illustration on the right),



$$P(X, Y, Z) = P(X) P(Y | X) P(Z | X, Y). \quad (1)$$

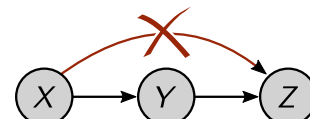
We then introduced the concept of *conditional independence* between two random variables X and Z given a third random variable Y , which is defined analogously to normal (i.e., unconditional) statistical independence as follows,

$$X \text{ and } Z \text{ are conditionally independent given } Y \Leftrightarrow P(X, Z | Y) = P(X | Y) P(Z | Y). \quad (2)$$

- (a) Show that conditional independence between X and Z given Y means that, once you know the value of Y , learning about the value of X would not provide any additional information about Z , i.e.,

$$X \text{ and } Z \text{ are cond. indep. given } Y \Leftrightarrow P(Z | X, Y) = P(Z | Y). \quad (3)$$

Remark: Eq. 3 implies that, if and only if X and Z are conditionally independent given Y , then the chain rule from Eq. 1 simplifies as follows (see also illustration on the right),



$$X \text{ and } Z \text{ are cond. indep. given } Y \Leftrightarrow P(X, Y, Z) = P(X) P(Y | X) P(Z | Y). \quad (4)$$

If three random variables X , Y , and Z satisfy the right-hand side of Eq. 4, then we say that they form a *Markov chain* $X \rightarrow Y \rightarrow Z$. A Markov chain can be interpreted as a memoryless process: if you want to generate a random sample from a Markov chain, then you can do so by *ancestral sampling*: you first generate some random sample $x \sim P(X)$, then you generate some $y \sim P(Y | X = x)$, and finally you generate some $z \sim P(Z | Y = y)$. Notice that, once you've generated the random sample y , you no longer need to keep x in memory because you won't need it for generating z . Later in this course, you will prove the important *data processing inequality*, which states that information about the initial random variable X can never increase along a Markov chain, i.e., $I_P(X; Z) \leq I_P(X; Y)$. The information processing inequality has far reaching consequences, e.g., on how information can propagate along a deep neural network.

(b) Show that conditional independence is neither a strictly stronger nor a strictly weaker property than normal (i.e., unconditional) independence. Thus,

- (i) show that two random variables X and Z can be statistically independent even if they are not conditionally independent given some third random variable Y ;

Hint: Consider our Simplified Game of Monopoly. You already showed in Problem 4.1 (b) that X_{red} and X_{blue} are statistically independent. Now show that X_{red} and X_{blue} are, however, *not* conditionally independent given X_{sum} .

- (ii) show that two random variables X and Z can be conditionally independent given some third random variable Y even if X and Z by themselves are not statistically independent.

Hint: Almost any Markov process will do. For example, you could consider a sequence of three independent coin tosses and let $C_i \in \{0, 1\}$ be the result of the i^{th} coin toss. Assume that the coin is bent (why is this necessary?), i.e., $P(C_i = 1) = \alpha$ and $P(C_i = 0) = 1 - \alpha$ with some $\alpha \neq \frac{1}{2}$. Then define $X := C_1$, $Y := (X + C_2) \bmod 2$, and $Z := (Y + C_3) \bmod 2$. Thus, X and Z are clearly conditionally independent given Y because they satisfy Eq. 3. Show by explicit calculation that X and Z , however, do *not* satisfy normal statistical independence.

Problem 5.2: Correlated Symbols in Various Probabilistic Model Architectures

In last week's lecture, we introduced various architectures for models of complicated probability distributions. We will make these architectures more concrete and discuss which compression method works best with which model architecture in upcoming lectures. But before we do this, let's analyze how capable each model architecture actually is. In particular, we analyze whether each of the proposed probabilistic models can capture *correlations* between symbols in a message, i.e., the fact that, in messages that appear in the real world, symbols are typically *not* statistically independent. All models below describe a message $\mathbf{X} = (X_1, X_2, \dots, X_k)$ where each symbol X_i , $i \in \{1, 2, \dots, k\}$ is modeled as a random variable with values from some discrete alphabet \mathfrak{X} .

The four parts (a)-(d) of this problem can be solved independently. So don't give up if you have trouble with some part.

- (a) **Fully factorized models:** before we look at more complicated model architectures below, let's consider the most trivial class of models which assume that all symbols X_i , $i \in \{1, 2, \dots, k\}$ are statistically independent. Such a model is often called "fully factorized" because the joint probability distribution $P(\mathbf{X})$ of the

message \mathbf{X} can be written as a product of the marginal distributions:

$$P_{\text{model}}(\mathbf{X}) = \prod_{i=1}^k P_{\text{model}}(X_i). \quad (5)$$

Here, we explicitly reinstated the subscript “model” because we now want to search for the best model, $P_{\text{model}}^*(\mathbf{X})$, in the form of Eq. 5 that approximates some data distribution $P_{\text{data}}(\mathbf{X})$, which is typically *not* fully factorized.

Thus, consider the cross entropy $H(P_{\text{data}}(\mathbf{X}), P_{\text{model}}(\mathbf{X}))$. Convince yourself that, for a model of the form of Eq. 5 (warning: but not for more general models),

$$H(P_{\text{data}}(\mathbf{X}), P_{\text{model}}(\mathbf{X})) = \sum_{i=1}^k H(P_{\text{data}}(X_i), P_{\text{model}}(X_i)) \quad (6)$$

where, following our usual notation, $P_{\text{data}}(X_i)$ is the marginal distribution of symbol X_i under P_{data} (i.e., the distribution that you obtain if you *marginalize* $P(\mathbf{X})$ over all X_j with $j \neq i$). Then argue that the right-hand side of Eq. 6 is minimized by setting $P_{\text{model}}^*(X_i) = P_{\text{data}}(X_i)$ for all i . (*Hint*: what is the cross entropy $H(P, P)$ of a distribution with itself, and why is it smaller or equal than any $H(P, Q)$ for all other distributions $Q \neq P$?)

Thus, within the class of fully factorized models, the best approximation $P_{\text{model}}^*(\mathbf{X})$ of an arbitrary distribution $P_{\text{data}}(\mathbf{X})$ is the product of the marginals, $P_{\text{model}}^*(\mathbf{X}) = \prod_{i=1}^k P_{\text{data}}(X_i)$. Convince yourself that, for this model, the cross entropy is the sum of the marginal entropies of each symbol under the data distribution,

$$H(P_{\text{data}}(\mathbf{X}), P_{\text{model}}^*(\mathbf{X})) = \sum_{i=1}^k H_{P_{\text{data}}}[X_i] \quad (\text{for optimal fully factorized model}). \quad (7)$$

- (b) **Markov Chains:** as discussed in the lecture, a Markov chain models the creation of a sequence of symbols X_1, X_2, \dots, X_k as a memoryless stochastic process, i.e.,

$$P(\mathbf{X}) = P(X_1) \prod_{i=2}^k P(X_i | X_{i-1}) \quad (8)$$

where, from here on, we drop the subscript “model” for simplicity.

- (i) Show that, although each symbol X_i is conditioned only on its immediately preceding symbol X_{i-1} (for $i > 1$) and not on any earlier symbols, a Markov chain can still model correlations between *any* symbols, even if they are not nearest neighbors. I.e., show that there exists a model of the form of Eq. 8 where all pairs of symbols X_i and X_j with $i \neq j$ are *not* statistically independent.

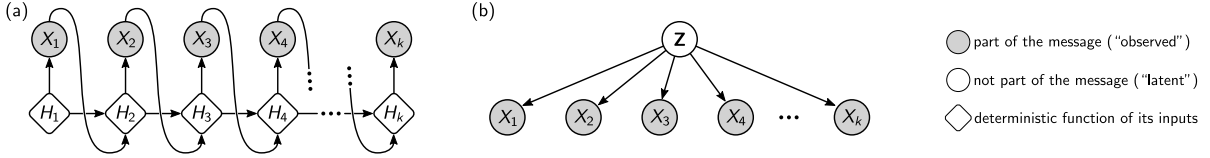


Figure 1: (a) autoregressive model, see Problem 5.2 (c); (b) latent variable model, see Problem 5.2 (d)

Hint: For example, you could consider the Markov chain over the alphabet $\mathcal{X} = \{0, 1\}$ with $P(X_1=0) = P(X_1=1) = \frac{1}{2}$ and

$$P(X_i = x_i | X_{i-1} = x_{i-1}) = \begin{cases} 0.99 & \text{for } x_i = x_{i-1}; \\ 0.01 & \text{for } x_i \neq x_{i-1}. \end{cases} \quad (9)$$

Then convince yourself (either by explicit calculation or simply by reasoning about what this Markov process models) that all marginal probabilities are $P(X_i=0) = P(X_i=1) = \frac{1}{2} \forall i$ by symmetry but that, e.g., the conditional probability $P(X_j=1 | X_i=1) > \frac{1}{2}$ for at least for some $i \neq j$ (it actually turns out that $P(X_j=1 | X_i=1) > 1 \forall j \geq i$ but this is more difficult to show).

- (ii) Now show that, although a Markov chain can model symbols that are not statistically independent, any two symbols X_i and X_l with $l \geq i + 2$ are *conditionally* independent given any X_j with $i < j < l$.

Hint: write out the joint probability of all symbols *up to* X_l as follows,

$$P(X_1, \dots, X_l) = \underbrace{\left(P(X_1) \prod_{\alpha=2}^i P(X_\alpha | X_{\alpha-1}) \right)}_{=P(X_1, \dots, X_i)} \underbrace{\left(\prod_{\alpha=i+1}^j P(X_\alpha | X_{\alpha-1}) \right)}_{=P(X_{i+1}, \dots, X_j | X_i)} \times \underbrace{\left(\prod_{\alpha=j+1}^l P(X_\alpha | X_{\alpha-1}) \right)}_{=P(X_{j+1}, \dots, X_l | X_j)}. \quad (10)$$

What do you get if you now marginalize both sides over all symbols except X_i , X_j , and X_l ? Compare the result to Eq. 4.

- (c) **Autoregressive models:** Figure 1 (a) illustrates an autoregressive model like the one you've used in Problem 3.2. The figure is a graphical representation of the following factorization of the joint probability distribution,

$$P(\mathbf{X}) = \prod_{i=1}^k P(X_i | H_i) \quad \text{with} \quad H_1 = \text{fixed}; \quad H_{i+1} = f(H_i, X_i) \quad (11)$$

where f is some deterministic function (e.g., a neural network). Show that autoregressive models are more powerful than Markov chains in that they can model probability distributions where two symbols X_i and X_l are *not* conditionally independent given some third symbol X_j with $i < j < l$.

Hint: For example, you could consider a toy autoregressive model over the alphabet $\mathfrak{X} = \{0, 1\}$ with $H_1 = 0$ and $H_{i+1} = f(H_i, X_i) = (H_i + X_i) \bmod 10$. Thus, the hidden state H_i counts how many “1” symbols have appeared before symbol X_i (modulo 10 so that the hidden states don’t grow out of bounds). Now you could make the probability of “1” symbols depend on H_i , e.g., by setting $P(X_i=1 | H_i) = \frac{H_i+1}{10}$ and $P(X_i=0 | H_i) = 1 - \frac{H_i+1}{10}$. Then, consider the first three symbols X_1, X_2 , and X_3 (the statement is also true for other triples of symbols, but the calculations are more tedious). Show by explicit calculation that

$$P(X_3=1 | X_1=1, X_2=1) \neq P(X_3=1 | X_2=1), \quad (12)$$

i.e., that even this simple toy model already violates the right-hand side of Eq. 3. The value of the left-hand side of Eq. 12 follows directly from the model definition but calculating the right-hand side takes a few more steps. Before you do these calculations, test your understanding by reasoning in words whether you expect the left-hand side of Eq. 12 to be smaller or larger than the right-hand side.

- (d) **Latent variable models:** Figure 1 (b) illustrates a latent variable model. You’ll learn how to use latent variable models for effective data compression with the so-called bits-back trick in the next lecture. But let’s first prove here that latent variable models can in fact capture correlations between symbols.

The illustration in Figure 1 (b) is a pictorial representation of the following factorization of a joint probability distribution over symbols $\mathbf{X} = (X_1, \dots, X_k)$ and a (usually multidimensional) so-called *latent* variable Z ,

$$P(\mathbf{X}, Z) = P(Z) \prod_{i=1}^k P(X_i | Z). \quad (13)$$

Here $P(Z)$ is called the “prior distribution” and $P(X_i | Z)$ is called the “likelihood”. At a first glance, the model architecture in Eq. 13 might look like it couldn’t possibly capture any correlations between different symbols X_i because the part of Eq. 13 that describes symbols is fully factorized (similar to the model in Eq. 5). However, this impression is deceptive because the symbols X_i are only *conditionally independent given the latent Z* . However, Z is not part of the message. The probabilistic model of the message is the *marginal* distribution of \mathbf{X} ,

$$P(\mathbf{X}) = \begin{cases} \sum_z P(\mathbf{X}, Z=z) & \text{for discrete } Z; \\ \int P(\mathbf{X}, Z=z) dz & \text{for continuous } Z. \end{cases} \quad (14)$$

Show that the marginal distribution in Eq. 14 can indeed describe correlations between symbols, i.e., a distribution of this form can model data sources where any two symbols X_i and X_l are *not* statistically independent, and are also *not* conditionally independent given any different third symbol X_j .

Hint: You could consider, e.g., a toy model over the alphabet $\mathfrak{X} = \{0, 1\}$ with $k = 3$, scalar $Z \in \{0, 1\}$, and with a likelihood $P(X_i | Z)$ that is the same for all i . Come up with some explicit probabilities for $P(Z = z)$ and $P(X_i = x_i | Z = z)$ for all $z, x_i \in \{0, 1\}$. Then show first that $P(X_1 = x_1, X_3 = x_3) \neq P(X_1 = x_1)P(X_3 = x_3)$ and finally that $P(X_3 = x_3 | X_1 = x_1, X_2 = x_2) \neq P(X_3 = x_3 | X_2 = x_2)$ in your model for some $x_1, x_2, x_3 \in \{0, 1\}$ of your choice. Try to explain your findings in words too: why does knowing the value of, e.g., X_1 influence the probability distribution over X_3 ?

No programming problem this week; next week's problem set will be mostly programming however. You'll improve our autoregressive compression method for natural language by replacing the Huffman coder with a range coder, and you'll implement and empirically analyze the bits-back trick for a toy latent variable model.