

# STATISTICAL METHODS FOR MACHINE LEARNING II

© by Xia, Tingfeng

2020 winter term

This work is licensed under a [Creative Commons “Attribution-NonCommercial-ShareAlike 4.0 International”](https://creativecommons.org/licenses/by-nc-sa/4.0/) license.



## Preface

This document is consist of notes from lectures and the online course notes that I, personally, find interesting/important. You can find the online course notes on the course website here: <https://probm1course.github.io/sta414/>

## Contents

<b>1</b>	<b>Lecture 2 - Introduction to Probabilistic Models</b>	<b>4</b>
1.1	Overview . . . . .	4
1.2	Probabilistic Perspective on ML . . . . .	4
1.2.1	(Example) Classification . . . . .	4
1.3	Observed vs Unobserved Random Variables . . . . .	5
1.3.1	Supervised Dataset . . . . .	5
1.3.2	Unsupervised Dataset . . . . .	5
1.3.3	Latent Variables . . . . .	5
1.4	Operations on Probabilistic Models . . . . .	6
1.5	Desiderata of Probabilistic Models . . . . .	6
1.5.1	Fully Dependent Factorization (Chain Rule) . . . . .	6
1.5.2	Assumptions (Independence) . . . . .	6
1.6	Likelihood Function . . . . .	7
1.7	Maximum Likelihood Estimation . . . . .	7
1.8	Sufficient Statistics . . . . .	7
1.8.1	Fisher-Neyman Factorization Theorem . . . . .	7
1.9	Exponential Family . . . . .	8
1.9.1	(Example) 1-D Gaussian . . . . .	8
<b>2</b>	<b>Lecture 3 - Directed Graphical Models</b>	<b>8</b>
2.1	Decision Theory (Utility Theory) . . . . .	8
2.2	Graphical Model Notation . . . . .	8
2.2.1	Chain Rule Expansion . . . . .	8

2.2.2	Graph Representation . . . . .	9
2.2.3	Conditional Independence . . . . .	9
2.2.4	Plates . . . . .	10
2.3	Directed Acyclic Graphical Models (DAGM) . . . . .	10
2.3.1	Independence Assumption on DAGMs . . . . .	10
2.3.2	(Example) Markov Chain . . . . .	11
2.4	Directed - Separation . . . . .	11
2.4.1	(Definition) D - Separation . . . . .	11
2.4.2	Float Rules Derivation . . . . .	11
2.4.3	The Bayes-Ball Algorithm . . . . .	13
2.4.4	Bayes Ball Rules . . . . .	13
2.5	Unobserved Variables . . . . .	14
2.5.1	Partially Unobserved Variables . . . . .	14
2.5.2	Latent Variables . . . . .	14
2.5.3	Mixture Models . . . . .	15
2.6	Examples . . . . .	15
2.6.1	Second-order Markov Chain . . . . .	15
2.6.2	Hidden Markov Models (HMMs) . . . . .	15
<b>3</b>	<b>Lecture 4 - Exact Inference</b>	<b>16</b>
3.1	Variable Elimination . . . . .	16
3.1.1	(Simple Example) Chain . . . . .	16
3.2	Sum-Product Inference . . . . .	17
3.2.1	SP - Inference in Directed Models . . . . .	17
3.2.2	SP - Inference in Undirected Models . . . . .	17
3.2.3	(Example) Directed Graph . . . . .	17
3.2.4	Complexity of VE . . . . .	18
<b>4</b>	<b>Lecture 5 - Message passing, Hidden Markov Models, and Sampling</b>	<b>19</b>
4.1	Message Passing . . . . .	19
4.1.1	Belief Propagation: Motivation and Definitions . . . . .	19
4.1.2	Belief Propagation Algorithm . . . . .	19
4.2	Inference in Hidden Markov Models . . . . .	19
4.3	The Forward - Backward Algorithm . . . . .	20
4.3.1	Forward Filtering Recursion . . . . .	21
4.3.2	Backward Filtering Recursion . . . . .	21
4.4	Sampling . . . . .	21
4.4.1	Ancestral Sampling . . . . .	21
4.4.2	Simple Monte Carlo . . . . .	22
<b>5</b>	<b>Lecture 7 - Stochastic Variational Inference (SVI/ADVI)</b>	<b>23</b>
5.1	Motivation - Approximating Posterior Inference . . . . .	23
5.2	The Kullback-Leibler Divergence . . . . .	23
5.2.1	Expectation Propagation . . . . .	24
5.3	The Evidence Lower BOund (ELBO) . . . . .	24
5.3.1	Definition and Equivalence Property . . . . .	24

5.3.2	Optimization goal of ELBO . . . . .	24
5.3.3	Pathwise Gradient . . . . .	25
5.4	Tutorial . . . . .	25
5.4.1	Derivation Using Jensen's Inequality . . . . .	25
5.4.2	Interpretations of ELBO . . . . .	26
<b>6</b>	<b>Lecture 8 - Sampling and Monte Carlo Methods</b>	<b>26</b>
6.1	Motivation: Sampling . . . . .	26
6.2	Lattice Discretization . . . . .	27
6.3	Importance Sampling . . . . .	27
6.4	Rejection Sampling . . . . .	29
6.5	Metropolis-Hasting Method: Example of MCMC . . . . .	30
<b>7</b>	<b>Lecture 9 - Amortized Inference and Variational Auto Encoders</b>	<b>30</b>
7.1	Motivations . . . . .	30
7.1.1	Efficient Use of (Accumulative) Partial Evidence . . . . .	30
7.2	Amortized Inference . . . . .	31
7.2.1	Simple Example (Parametrizing Gaussian) . . . . .	31
7.2.2	Amortized Inference - Algorithmic Procedure . . . . .	31
7.2.3	Optimizing Model Parameters . . . . .	31
7.3	Variational Auto Encoders (VAEs) . . . . .	32
7.3.1	Autoencoders: Definition and Problems . . . . .	32
7.3.2	Problems with Deterministic Autoencoders . . . . .	32
7.3.3	Variational Auto Encoders . . . . .	33
7.3.4	Utilization of Latent Variable Models . . . . .	33

# 1 Lecture 2 - Introduction to Probabilistic Models

## 1.1 Overview

We have a random vector in the form  $X = (X_1, \dots, X_m)$  which can be *either observed or unobserved*. To approach this in a generative way, we make the so called generative assumption. which is that  $X \sim P_{true}(X)$ , i.e. there is some true distribution that is behind the scene and our data is from such distribution.

**Goal** Model a parametric joint distribution  $P_\theta(X)$  by learning the parameters. The learning here means we want to find a/the “close”/“best” estimation to our parameter  $\theta$ . In this course we will investigate the following three problems,

- How to specify the joint,  $P_\theta(X)$ ?
- What does “best”/“close” mean? In some sense we want to find  $P_\theta \approx P_{true}$ , however  $P_{true}$  might also be unknown.
- How to find the best  $\theta$ ? In this course we will generally rely on gradient methods, so  $\nabla_{\theta} \dots$

## 1.2 Probabilistic Perspective on ML

With this perspective, we can think about common machine learning tasks differently, where random variables represent:

- $X$ : (high dimensional) input data
- $C$ : discrete label
- $Y$ : continuous target

If we assume our knowledge of the joint of the above three, i.e. we know  $P(X, C, Y)$ , then we can write our familiar tasks in the following way

- **Regression:**

$$p(Y|X) = \frac{p(X, Y)}{P(X)} = \frac{p(X, Y)}{\int p(X, Y) dY} \quad (1.1)$$

- **Classification/Clustering:**

$$p(C|X) = \frac{p(X, C)}{\sum_{C'} p(X, C')} \quad (1.2)$$

### 1.2.1 (Example) Classification

Suppose we have data of the form  $\mathcal{D} = \{(x, c)_i\}_i$ . We assume that they came from a certain true distribution, i.e.  $\{(x, c)_i\}_i \sim p(X, C)$ . Then, the ultimate goal of the ML problem is converted into finding  $p(C|X)$ . Using Bayes Rule of total probability, we can expand the distribution of interest into

$$p(C|X) = \frac{p(X, C)}{P(X)} = \frac{p(X, C)}{\sum_{C'} p(X, C')} \quad (1.3)$$

**Output Heuristics** After we acquire  $p(C|X)$  as above, we are one step away from our goal of output the actual prediction  $c^*$ . There are three ways that we can do this, namely

- **MLE Estimate** is the most intuitive one, we simply choose

$$c^* = \arg \max_c p(C = c|X) \quad (1.4)$$

- **Sample Learnt Dist** is another approach which produces non-deterministic results, i.e. we sample  $c^* \sim p(C|X)$ .
- **Combined** is usually a safe way of doing this. We output

$$(c^*, p(C = c^*|X)) \iff (\text{result, how sure?}) \quad (1.5)$$

As an example, we have a ML algorithm drives a car. In this case, we might want to make decision only when the machine learning model has a certain level of confidence.

### 1.3 Observed vs Unobserved Random Variables

#### 1.3.1 Supervised Dataset

$$\{x_i, c_i\}_{i=1}^N \sim p(X, C) \quad (1.6)$$

In such case, the class labels are observed and finding the conditional distribution  $p(C|X)$  satisfies the supervised classification problem.

#### 1.3.2 Unsupervised Dataset

$$\{x_i\}_{i=1}^N \sim p(X, C) \quad (1.7)$$

Still under the generative assumption, where we assume that there is some underlying distribution for our dataset. Further, we assume that the distribution of data is related to the class labels for the data points even though the class labels are never observed. **A common way to refer to an unobserved discrete class label is “cluster”**. However, in this case, our final goal of classification is still  $p(C|X)$ <sup>1</sup>.

#### 1.3.3 Latent Variables

Further, like clusters, introducing assumptions about unobserved variables is a powerful modelling tool. We will make use of this by modelling variables which are never observed in the dataset, called latent or hidden variables. By introducing and modelling latent variables, we will be able to naturally describe and capture abstract features of our input data.

---

<sup>1</sup>Might be helpful to think of Gaussian Mixture Models

## 1.4 Operations on Probabilistic Models

- **Generate Data:** Sample from the model.
- **Estimate Likelihood:** When all variables are either observed or marginalized, we produce the result which is a single real number that describes the ‘probability’ of the all variables taking on those specific values.
- **Inference:** Compute the expected value of some variables given others which are either observed or marginalized.
- **Learning:** Set the parameters of the joint distribution given some observed data to maximize the probability of the observed data.

## 1.5 Desiderata of Probabilistic Models

We have two desires for the joint distribution to learn, namely

- The marginal and conditional distribution can be computed efficiently
- The representation of the joint distribution should be compact. This is especially important when we are dealing with joint distributions over many variables.

In general, total joint distribution are too large to specify and would require an insane amount of data to fit even we wanted to. Thus, we need modelling assumptions.

### 1.5.1 Fully Dependent Factorization (Chain Rule)

Suppose we have sample space defined such that  $|T| = 2$ ,  $|W| = 3$ , and  $|M| = 4$ . Then the total joint distribution could be expanding using the chain rule as (*note: not unique*)

$$P_{\theta}(T, W, M) = P(T)P(W|T)P(M|T, W) \quad (1.8)$$

which requires<sup>2</sup>  $|\theta| = (2 - 1) + (3 - 1) \times 2 + (4 - 1) \times 2 \times 3 = 23$  parameters in total to specify.

### 1.5.2 Assumptions (Independence)

Introducing assumptions results in

- a less expressive model
- $|\theta|$  (usually, much) smaller

which can be bad sometimes (since the model is less expressive) and is a trade-off that we as modellers have to deal with.

**(Example) Fully Independent** We assume that  $T \perp W \perp M$ , then by definition we know, for example,  $P(W|T) = P(W)$ . Then

$$P_{\theta}(T, W, M) = P(T)P(W|T)P(M|T, W) \quad (1.9)$$

$$= P(T)P(W)P(M) \quad (1.10)$$

which requires only  $|\theta| = (2 - 1) + (3 - 1) + (4 - 1) = 6$  to fit.

---

<sup>2</sup>Here the bars mean “cardinality” rather than vector norm

## 1.6 Likelihood Function

For some observed data  $X$ , the likelihood describes the likeliness of the data under then distribution with parameter  $\theta$ .

$$L(\theta) = p(X|\theta) \quad (1.11)$$

In general, we prefer to deal with the log likelihood function, which is defined as

$$\ell(\theta; X) = \log L(\theta) = \log p(X|\theta) \quad (1.12)$$

## 1.7 Maximum Likelihood Estimation

The idea is to find

$$\hat{\theta}_{MLE} := \arg \max_{\theta} \ell(\theta; \mathcal{D}) \quad (1.13)$$

In the case of i.i.d, we can re-write as

$$\hat{\theta}_{MLE} = \arg \max_{\theta} \ell(\theta; \mathcal{D}) \quad (1.14)$$

$$= \arg \max_{\theta} \log \prod_m p(x^{(m)}|\theta) \quad (1.15)$$

$$= \arg \max_{\theta} \sum_m \log p(x^{(m)}|\theta) \quad (1.16)$$

## 1.8 Sufficient Statistics

**(Definition) Statistic** A statistic is a possibly vector valued *deterministic* function of a set of random variables.

**(Definition) Sufficient Statistic** is a statistic that conveys exactly the same information about the data generating process that created the data as the entire data itself.<sup>3</sup> In formal language, Sufficient Statistic ( $T(X)$ ) for  $X$  could be defined as

$$T(X^{(1)}) = T(X^{(2)}) \implies L(\theta; X^{(1)}) = L(\theta; X^{(2)}), \quad \forall \theta \quad (1.17)$$

alternatively, we can define it as

$$P(\theta|T(X)) = P(\theta|X); \quad \text{i.e. data doesn't give further info} \quad (1.18)$$

### 1.8.1 Fisher-Neyman Factorization Theorem

If the probability function is  $f_{\theta}(x)$ , then  $T$  is a sufficient statistic for  $\theta$  if and only if non-negative functions  $g$  and  $h$  can be found such that

$$P(\theta|T(X)) = h(x, T(x))g(T(x), \theta) \quad (1.19)$$

---

<sup>3</sup>Could also be interpreted as “summarize the data with respect to the likelihood”

## 1.9 Exponential Family

Factorizes as

$$p(x|\eta) = h(x) \exp \{ \eta' T(x) - g(\eta) \} \quad (1.20)$$

$$= h(x) g(\eta) \exp \{ \eta' T(x) \} \quad (1.21)$$

### 1.9.1 (Example) 1-D Gaussian

$$p(x|\theta) = \mathcal{N}(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left( -\frac{1}{2\sigma^2} (x - \mu)^2 \right) \quad (1.22)$$

$$= \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left( -\frac{1}{2\sigma^2} (x^2 - 2x\mu + \mu^2) \right) \quad (1.23)$$

$$= \underbrace{\frac{1}{\sqrt{2\pi\sigma}}}_{h(x)} \underbrace{\exp \left( \frac{-\mu^2}{2\sigma^2} \right)}_{g(\eta)} \exp \left( \underbrace{\left[ \frac{\mu}{\sigma^2} \quad \frac{-1}{2\sigma^2} \right]}_{\eta} \underbrace{\begin{bmatrix} x \\ x^2 \end{bmatrix}}_{T(X)} \right) \quad (1.24)$$

## 2 Lecture 3 - Directed Graphical Models

### 2.1 Decision Theory (Utility Theory)

Let  $a$  denote action, and  $a^*$  be the optimal one. Also use  $s$  to denote state and  $V(\cdot)$  be the value function. We have, in general

$$a^* = \arg \min_a / \arg \max_a \underbrace{\mathbb{E}_{p(s|a, \text{knowledge})} [V(s)]}_{u(a) \triangleq \text{utility of action } a} \quad (2.1)$$

Often, it is very hard to find a/the good/best value function that quantifies everything using a numerical value.

### 2.2 Graphical Model Notation

#### 2.2.1 Chain Rule Expansion

Given any joint probability of  $N$  random variables, we can expand it as follows

$$p(x_1, \dots, x_N) = p(x_1) p(x_2|x_1) p(x_3|x_2, x_1) \dots p(x_N|x_{N-1:1}) \quad (2.2)$$

Formally speaking, in the case of two random variables would simply

$$p(x, y) = p(x|y)p(y) \quad (2.3)$$

and the general case for  $N$  random variables could be written as <sup>4</sup>

$$p \left( \bigcap_{i=1}^N x_i \right) = \prod_{j=1}^N p \left( x_j \left| \bigcap_{k=1}^{j-1} x_k \right. \right) \quad (2.4)$$

<sup>4</sup>Note: when  $k = 1$ ,  $p(x_k | \bigcap_{j=1}^{k-1} x_j) = p(x_1)$

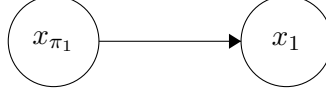


### 2.2.2 Graph Representation

**(Example) Grouping Variables** Consider the model

$$p(x_i, x_{\pi_i}) = p(x_{\pi_i}) p(x_i | x_{\pi_i}) \quad (2.5)$$

which we can use the following graph to represent:

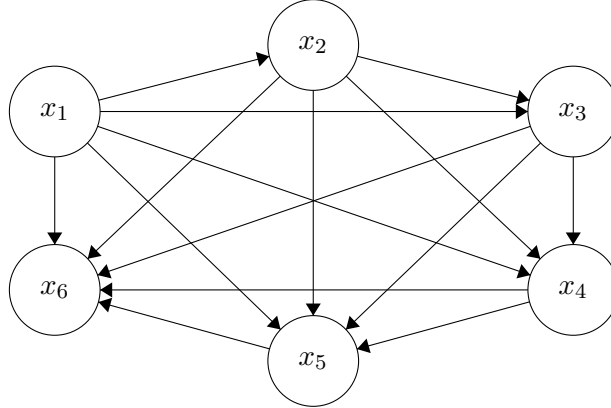


where

- **nodes** represent random variables
- **arrows** mean “conditioned on”, e.g. “ $x_i$  is conditioned on  $x_{\pi_i}$ ”

Notice that we can always group the variables together into one bigger variable, so in this example,  $x_{\pi_i}$  might represent a group of variables instead of just one.

**(Example) Fully Dependent 6 Nodes** The total expansion of  $p(x_{1:6})$  could be represented as



This is the resultant graphical model if we make **absolutely no assumption** on the independence. Notice that such model grows exponentially in complexity with respect to the number of parameters considered. We say such model “scales poorly”.

### 2.2.3 Conditional Independence

**(Definition) Conditional Independence** Let  $X$  be the set of nodes in our graph (the random variables of our model), then two sets of variables  $X_A, X_B$  are said to be conditionally independent given a third set of variables  $X_C$  if and only if either

$$p(X_A, X_B | X_C) = p(X_A | X_C) p(X_B | X_C) \quad (2.6)$$

or (**#! Important!**)

$$p(X_A | X_B, X_C) = p(X_A | X_C) \quad (2.7)$$

$$\iff p(X_B | X_A, X_C) = p(X_B | X_C) \quad (2.8)$$

and we denote the relation as  $(X_A \perp X_B | X_C)$

### 2.2.4 Plates

In Bayesian methods, we treat parameters as random variables and hence we would like to include them in our graphical model. However, adding a node for each observation is quite cumbersome and thus we introduce plates, which denote replication of random variables.

**Nested Plates** Plates could be nested, in which case their arrows get duplicated also, *according to the rule:* draw an arrow from every copy of the source node to every copy of the destination node.

**Crossing Plates** Plates can also cross (intersect), in which case the nodes at the intersection have multiple indices and get duplicated a number of times equal to the product of the duplication numbers on all the plates containing them.

## 2.3 Directed Acyclic Graphical Models (DAGM)

A directed acyclic graphical model over  $N$  random variables look like

$$p(x_{1:N}) = \prod_i^N p(x_i | x_{\pi_i}) \quad (2.9)$$

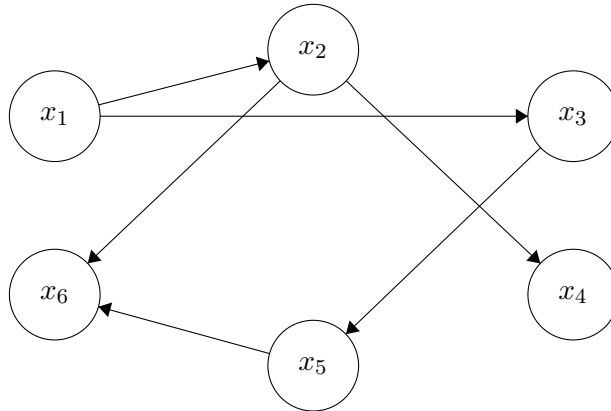
where  $x_i$  is a random variable and  $x_{\pi_i}$  denotes the parents of the node (which could be an empty set). This notion is more general than the fully dependence model that looked at above. Notice that here each node is only dependent on its parents rather than all other nodes. Thus, the complexity of such model reduces to exponential in fan-in of each node, instead of the total  $N$ .

### 2.3.1 Independence Assumption on DAGMs

Then, we have the independence relationship of<sup>5</sup>  $x_i \perp x_{Ancestor(\pi_i)} | x_{\pi_i}$  which expands into

$$p(x_{1,\dots,6}) = p(x_1) p(x_2|x_1) p(x_3|x_1) p(x_4|x_2) p(x_5|x_3) p(x_6|x_2, x_5) \quad (2.10)$$

with the following respective graph




---

<sup>5</sup>Requires topological ordering, to be added later.

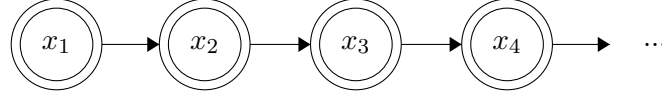
As we can see, the introduction of the assumption greatly reduced the complexity of the model.

### 2.3.2 (Example) Markov Chain

The following example has independence relationships that satisfies the Markov Property.

$$p(x) = p(x_1)p(x_2|x_1)p(x_3|x_2)p(x_4|x_3)\dots \quad (2.11)$$

and could be represented, in graphical model, as



## 2.4 Directed - Separation

### 2.4.1 (Definition) D - Separation

Directed-separation is a notion of connectedness in DAGs in which two (sets of) variables may or may not be connected conditioned on a third (set of) variable(s). D-connection implies conditional dependence and d-separation implies conditional independence.

### 2.4.2 Float Rules Derivation

**Notation** The double circles (or shaded circles) represent the notion of “conditioned-on”.

**Chain**



In this case, we are interested in knowing whether or not

$$(X \perp Z) | Y \quad (2.12)$$

From the arrows between the nodes, we know that

$$P(X, Y, Z) = P(X)P(Y|X)P(Z|Y) \quad (2.13)$$

Then, we have

$$P(X, Z|Y) = \frac{P(X, Y, Z)}{P(Y)} \quad (2.14)$$

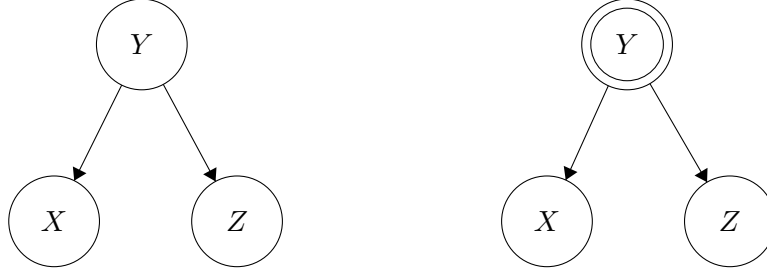
$$= \frac{P(X)P(Y|X)P(Z|Y)}{P(Y)} \quad (2.15)$$

$$= \frac{P(X, Y)P(Z|Y)}{P(Y)} \quad (2.16)$$

$$= P(X|Y)P(Z|Y) \quad (2.17)$$

and this completes the proof. ■

### Common Clause



As previous, we are interested in knowing whether or not  $(X \perp Z)|Y$ . From the graph, we can know that

$$P(X, Y, Z) = P(Y)P(X|Y)P(Z|Y) \quad (2.18)$$

Then, we have

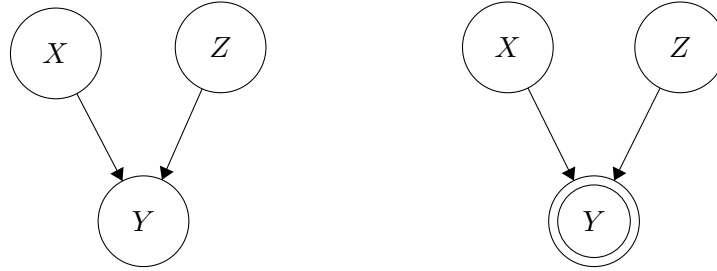
$$P(X, Z|Y) = \frac{P(X, Y, Z)}{P(Y)} \quad (2.19)$$

$$= \frac{P(Y)P(X|Y)P(Z|Y)}{P(Y)} \quad (2.20)$$

$$= P(X|Y)P(Z|Y) \quad (2.21)$$

and this completes the proof. ■

### Explaining Away (Berkson's Paradox)



Notice from the graph that we have  $P(X, Y, Z) = P(X)P(Z)P(Y|X, Z)$ . First, I will prove that, in this case,  $(X \not\perp Z)|Y$ .

$$P(Z|X, Y) = \frac{P(X)P(Z)P(Y|X, Z)}{P(X)P(Y|X)} \quad (2.22)$$

$$= \frac{P(Z)P(Y|X, Z)}{P(Y|X)} \quad (2.23)$$

$$\neq P(Z|Y) \quad (2.24)$$

For marginal independence, i.e.  $(X|Z)$ , we want to show that  $P(X, Z) = P(X)P(Z)$ .

$$P(X, Z) = \sum_{Y'} P(X, Y', Z) \quad (2.25)$$

$$= \sum_{Y'} P(X)P(Z)P(Y'|X, Z) \quad (2.26)$$

$$= P(X)P(Z) \sum_{Y'} P(Y'|X, Z) \quad (2.27)$$

$$= P(X)P(Z) \quad (2.28)$$

■

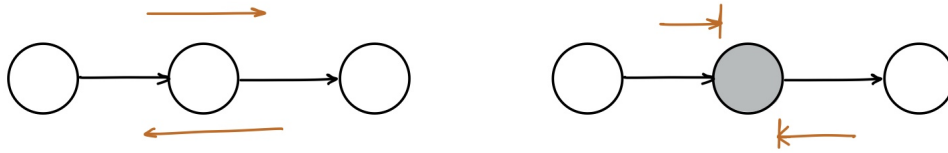
### 2.4.3 The Bayes-Ball Algorithm

To check if  $x_A \perp x_B | x_C$ , we will need

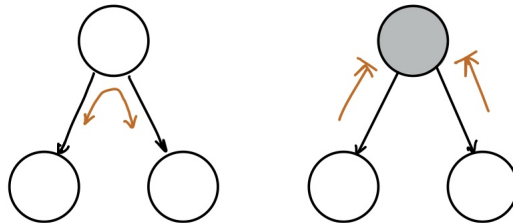
- Shade all nodes that were conditioned on, i.e. all nodes of  $x_C$
- Place balls at each node in  $x_A$  (or  $x_B$ )
- Let the balls bounce around according to rules that are yet to be states below
  - If any of the balls reach any of the nodes in  $x_B$  from  $x_A$  (or reach  $x_A$  from  $x_B$ ) then we declare  $x_A \not\perp x_B | x_C$
  - Otherwise,  $x_A \perp x_B | x_C$

### 2.4.4 Bayes Ball Rules

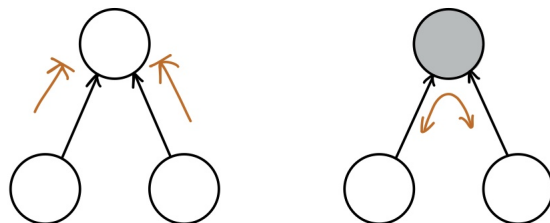
**Chain**



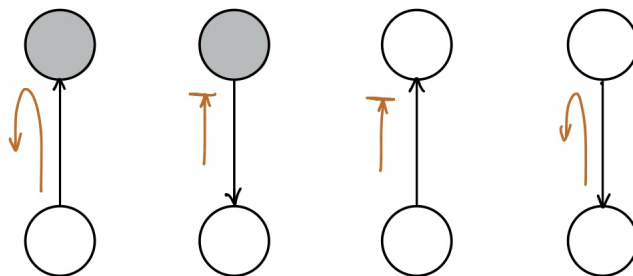
**Common Cause**



## Explain Away



## Linear



## 2.5 Unobserved Variables

Certain variables in our models may be unobserved, either some of the time or always, at training time or at test time. Graphically, we use shading to indicate observation.

### 2.5.1 Partially Unobserved Variables

<sup>6</sup> If variables are occasionally unobserved then they are missing data, e.g., undefined inputs, missing class labels, erroneous target values. In this case, we can still model the joint distribution, but we marginalize the missing values

$$\ell(\theta; \mathcal{D}) = \sum_{\text{complete}} \log p(x^c, y^c | \theta) + \sum_{\text{missing}} \log p(x^m | \theta) \quad (2.29)$$

$$= \sum_{\text{complete}} \log p(x^c, y^c | \theta) + \sum_{\text{missing}} \log \sum_y p(x^m, y | \theta) \quad (2.30)$$

### 2.5.2 Latent Variables

Above we discussed the case where some data are non-deterministically unobserved. Latent variables refers to those that **never observed**. The handling of the latent variables depends on where it appears in our model,

- If we never condition on it when computing the probability of the variables we do observe, then we can just forget about it and integrate it out. For example, given  $y$ ,  $x$  we fit the model

$$p(z, y | x) = p(z | y) p(y | x, w) p(w) \quad (2.31)$$

---

<sup>6</sup>An concrete example would be hospital data, where typically a large proportion of the data is missing.

- If  $z$  is not a leaf node, marginalizing over it will induce dependencies between its children. For example, given  $y, x$  we can fit the model

$$p(y|x) = \sum_z p(y|x, z)p(z) \quad (2.32)$$

### 2.5.3 Mixture Models

In this case, we are looking at data that has no input on class information. We can sum the labels out,

$$p(x|\theta) = \sum_{k=1}^K p(z = k|\theta_z) p(x|z = k, \theta_k) \quad (2.33)$$

where bayes rule comes in handy for calculating the posterior (class responsibilities) of the mixture component given some data, i.e.,

$$p(z = k|x, \theta_z) = \frac{p(z = k|\theta_z) p_k(x|\theta_k)}{\sum_j p(z = j|\theta_z) p_j(x|\theta_j)} \quad (2.34)$$

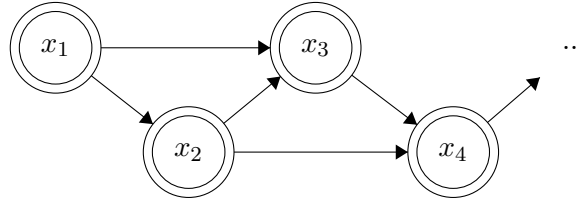
## 2.6 Examples

### 2.6.1 Second-order Markov Chain

Consider the model

$$p(\mathbf{x}_{1:T}) = p(x_1, x_2) p(x_3|x_1, x_2) p(x_4|x_2, x_3) \cdots = p(x_1, x_2) \prod_{t=3}^T p(x_t|x_{t-1}, x_{t-2}) \quad (2.35)$$

which has the graphical representation (double circles mean “observed” here)

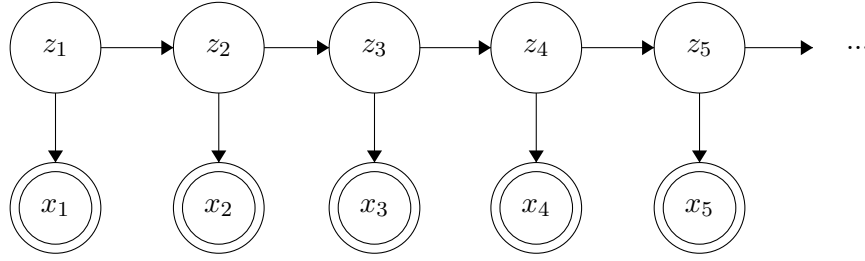


we notice that this model essentially assumes “the present depends on the past only through the current state as well as the last one.”

### 2.6.2 Hidden Markov Models (HMMs)

HMM is a statistical model in which a system being modelled is assumed to be a Markov Process<sup>7</sup> with un-observed states. Here is the graphical model, where double circles mean “observed” here:

<sup>7</sup>In continuous time, the Markov Chain model is known as Markov Process



where

- $z_t$  are hidden states taking on one of  $K$  discrete values
- $x_t$  are observed variables taking on values in any space.

The above graph factorizes into

$$p(X_{1:T}, Z_{1:T}) = p(Z_{1:T}) p(X_{1:T}|Z_{1:T}) = p(Z_1) \prod_{t=2}^T p(Z_t|Z_{t-1}) \prod_{t=1}^T p(X_t|Z_t) \quad (2.36)$$

## 3 Lecture 4 - Exact Inference

### 3.1 Variable Elimination

#### 3.1.1 (Simple Example) Chain

The example that we will consider is the simple chain

$$A \rightarrow B \rightarrow C \rightarrow D \quad (3.1)$$

where we want to compute  $P(D)$ , with no observation for other variables. We have

$$X_F = \{D\}, X_E = \{\}, X_R = \{A, B, C\} \quad (3.2)$$

The graphical model gives the factorization of

$$p(A, B, C, D) = p(A)p(B|A)p(C|B)p(D|C) \quad (3.3)$$

thus, if we want to find  $p(D)$  we can marginalize over all other variables, i.e.

$$p(D) = \sum_{A,B,C} p(A, B, C, D) \quad (3.4)$$

$$= \sum_C \sum_B \sum_A p(A)p(B|A)p(C|B)p(D|C) \quad (3.5)$$

If we marginalize the above the naïve way, then it cost exponentially  $\mathcal{O}(k^n)$ . Thus, we need to find a **elimination ordering** that gives us a smaller runtime complexity. We can reduce the complexity by first computing terms that appear across the other marginalization sums,

$$\phi(D) = \sum_C p(D|C) \sum_B p(C|B) \sum_A p(A)p(B|A) \quad (3.6)$$

$$= \sum_C p(D|C) \sum_B p(C|B)\phi(B) \quad (3.7)$$

$$= \sum_C p(D|C)\phi(C) \quad (3.8)$$

No observations for other variables means that we have no 'evidence' for other variables.

$X_F$  is the set of variable that we are interested in,  $X_E$  is evidence, and  $X_R$  is the set of extraneous variables, i.e. what we marginalize out

Here  $\phi$  means factor, an non-normalized "pmf"



## 3.2 Sum-Product Inference

Let  $X = Z \cup Y \wedge Z \cap Y = \emptyset$  be a set of random variables. . Then, to compute  $P(Y)$  for a directed and undirected model could be achieved by the sum product inference algorithm

$$\tau(Y) = \sum_z \prod_{\phi \in \Phi} \phi(\text{Scope}(\phi) \cap Z, \text{Scope}(\phi) \cap Y) \quad \forall Y \quad (3.9)$$

Here  $Y$  is the set that we want to do inference on and thus we are marginalizing over everything in  $Z$

where the scope for each  $\phi$  simply means the set of all random variables that have appeared in that specific  $\phi$ , and  $\Phi$  denotes the set of all  $\phi$ 's.

how to generalize to undirected graph?

### 3.2.1 SP - Inference in Directed Models

In a directed model, the  $\Phi$  is given by the conditional probability distributions for all variables, i.e.

$$\Phi = \{\phi_{x_i}\}_{i=1}^N = \{p(x_i | \text{parents}(x_i))\}_{i=1}^N \quad (3.10)$$

The resulting term  $\tau(Y)$  will be automatically normalized.

### 3.2.2 SP - Inference in Undirected Models

For undirected models,  $\Phi$  is given by the set of un-normalized potentials, and we **must** normalize the resulting  $\tau(Y)$  by  $\sum_Y \tau(y)$

### 3.2.3 (Example) Directed Graph

We have the following factorization of the joint distribution

$$p(C, D, I, G, S, L, H, J) = p(C)p(D|C)p(I)p(G|D, I)p(L|G)P(S|I)p(J|S, L)p(H|J, G) \quad (3.11)$$

and we will need the following potentials

$$\Phi = \{\phi(C), \phi(C, D), \phi(I), \phi(G, D, I), \phi(L, G), \phi(S, I), \phi(J, S, L), \phi(H, J, G)\} \quad (3.12)$$

consider the problem where we want to infer  $P(J)$  with elimination ordering  $\prec_{\{C,D,I,H,G,S,L\}}$ .<sup>8</sup> Then,

$$p(J) = \sum_L \sum_S \phi(J, L, S) \sum_G \phi(L, G) \sum_H \phi(H, G, J) \sum_I \phi(S, I) \phi(I) \sum_D \phi(G, D, I) \sum_C \phi(C) \phi(C, D) \quad (3.13)$$

$$= \dots \underbrace{\sum_C \phi(C) \phi(C, D)}_{\tau(D)} \quad (3.14)$$

$$= \dots \underbrace{\sum_D \phi(G, D, I) \tau(D)}_{\tau(G, I)} \quad (3.15)$$

$$= \dots \underbrace{\sum_I \phi(S, I) \phi(I) \tau(G, I)}_{\tau(S, G)} \quad (3.16)$$

$$= \dots \tau(S, G) \underbrace{\sum_H \phi(H, G, J)}_{\tau(S, G) \tau(G, J)} \quad // \text{ Note that } \tau(S, G) \text{ doesn't contain } H \quad (3.17)$$

$$= \dots \underbrace{\sum_G \phi(L, G) \tau(S, G) \tau(G, J)}_{\tau(J, L, S)} \quad (3.18)$$

$$= \dots \underbrace{\sum_S \phi(J, L, S) \tau(J, L, S)}_{\tau(J, L)} \quad (3.19)$$

$$= \dots \underbrace{\sum_L \tau(J, L)}_{\tau(J)} \quad (3.20)$$

$$= \tau(J) \quad (3.21)$$

### 3.2.4 Complexity of VE

The complexity of VE is

$$\mathcal{O}(mk^{N_{\max}}) \quad (3.22)$$

where

- $m = |\Phi|$  is the number of initial factors, i.e. the number of terms in the original factorization of the joint probability (given the graphical model)
- $k$  is the number of states each random variable takes (assumed to be equal here)
- $N_i$  is the number of random variables inside each sum  $\sum_i$  *at the time of marginalization*

---

<sup>8</sup>Note that what comes first in the ordering is the what we want to sum over in the inner most summation

- $N_{max} \triangleq \max_i N_i$  is the number of random variables inside the largest sum. (Inner sums are counted as one term all together)

## 4 Lecture 5 - Message passing, Hidden Markov Models, and Sampling

### 4.1 Message Passing

#### 4.1.1 Belief Propagation: Motivation and Definitions

Our goal is to compute the marginal of every variable in graph  $p(x_i), \forall x_i \in X$ . Notice that in a tree, we have

$$P(X_{1:n}) = \frac{1}{Z} \prod_{k \in \{1, \dots, n\}} \phi(x_k) \prod_{(i,j) \in T} \phi_{i,j}(x_i, x_j) \quad (4.1)$$

Below,  $T$  means the set of edges in the tree

and thus, if we want to compute  $P(X_1)$  we can marginalize out all other variables

$$P(X_1) = \sum_{x_2, \dots, x_n} P(X_{1:n}) = \sum_{x_2, \dots, x_n} \frac{1}{Z} \prod_k \phi(x_k) \prod_{(i,j) \in T} \phi_{i,j}(x_i, x_j) \quad (4.2)$$

$$\propto \sum_{x_2, \dots, x_n} \prod_k \phi(x_k) \prod_{(i,j) \in T} \phi_{i,j}(x_i, x_j) \quad (4.3)$$

Now that we have the goal, we can group what we want to compute into a more structured representation, where we define **message-passing** from variable  $j$  to  $i \in N(j)$  as

$$m_{j \rightarrow i}(x_i) = \sum_{x_j} \phi_j(x_j) \phi_{ij}(x_i, x_j) \prod_{k \in N(j) \neq i} m_{k \rightarrow j}(x_j) \quad (4.4)$$

#### 4.1.2 Belief Propagation Algorithm

1. Choose an root  $r$  arbitrarily,
2. Pass messages from leaves to  $r$
3. Pass messages from  $r$  to leaves
4. Compute

$$p(x_i) \propto \phi_i(x_i) \prod_{j \in N(i)} m_{j \rightarrow i}(x_i), \forall_i \quad (4.5)$$

Step 2 and 3 here so that we will compute all the  $m_{i \rightarrow j}, \forall (i, j), (j, i) \in T$ , after which we can use them to calculate the marginals using formula in step 4

### 4.2 Inference in Hidden Markov Models

There are the following four kinds of inference that we can do on a HMM, namely

- **Filtering:** Compute the belief state  $p(z_t | \mathbf{x}_{1:t})$  online. Notice that we call this inference operation filtering because this produces a result smoother than simply computing  $p(z_t | \mathbf{x}_t)$ .

Notice that the idea is that we want to calculate **all** of the marginals

- **Smoothing:** Compute  $p(z_t|\mathbf{x}_{1:T})$  offline. Notice that this computation happens after all  $\mathbf{x}_{1:T}$  observations have been collected and hence is offline.<sup>9</sup>
- **Fixed lag smoothing:** Compute  $p(z_{t-\ell}|\mathbf{x}_{1:t})$ , where  $\ell > 0$  is called the lag. Essentially, we want to do inference on the hidden state  $z_\alpha$  and we wait until we observe all of  $\mathbf{x}_{1:\alpha+\ell}$  before we carry out the computation.
- **Prediction:** In this case, we want to compute  $p(z_{t+h}|\mathbf{x}_{1:t})$  for some  $h > 0$ . ( $h$  is called the prediction horizon) Let's first take a look at the example where  $h = 2$ , then

Definition in KPM is different from course note

$$p(z_{t+2}|\mathbf{x}_{1:t}) = \sum_{z_{t+1}} \sum_{z_t} p(z_t, z_{t+1}, z_{t+2}|\mathbf{x}_{1:t}) \quad (4.6)$$

$$= \sum_{z_{t+1}} \sum_{z_t} p(z_t|\mathbf{x}_{1:t}) p(z_{t+1}|z_t, \mathbf{x}_{1:t}) p(z_{t+2}|z_{t+1}, z_t, \mathbf{x}_{1:t}) \quad (4.7)$$

$$= \sum_{z_{t+1}} \sum_{z_t} p(z_{t+2}|z_{t+1}) p(z_{t+1}|z_t) p(z_t|\mathbf{x}_{1:t}) \quad (4.8)$$

where the last step of simplification could be easily justified using Bayes Ball rules. Above, we computed the prediction about the future hidden states, and it can be converted into prediction about the future observations using

$$p(\mathbf{x}_{t+h}|\mathbf{x}_{1:t}) = \sum_{z_{t+h}} p(\mathbf{x}_{t+h}|z_{t+h}) p(z_{t+h}|\mathbf{x}_{1:t}) \quad (4.9)$$

which is called *Posterior Predictive Density*.

### 4.3 The Forward - Backward Algorithm

Suppose that we want to find  $p(z_t|\mathbf{x}_{1:T}), \forall t$ , then it can be broken down into two pieces

- **Forward Filtering:** Compute  $p(z_t, x_{1:t}) \quad \forall t$ , which effectively is the probability of the hidden state  $z_t$  given all past observations, up to including  $x_t$ .
- **Backward Filtering:** Compute  $p(x_{1+t:T}|z_t) \quad \forall t$ , which computes the probability of *all* future observations from  $x_{t+1}$  up to including  $x_T$  given the current hidden state, i.e.  $z_t$ .

Then, the complete smoothing  $p(z_t|\mathbf{x}_{1:T})$  could be computed as

$$p(z_t|\mathbf{x}_{1:T}) = p(x_{1:T}) p(z_t, x_{1:T}) \quad (4.10)$$

$$\propto p(z_t, x_{1:T}) \quad (4.11)$$

$$= p(z_t, x_{1:t}) p(x_{t+1:T}|z_t, x_{1:t}) \quad (4.12)$$

$$= p(z_t, x_{1:t}) p(x_{t+1:T}|z_t) \quad // \text{ By Bayes Ball} \quad (4.13)$$

$$= (\text{Forward Recursion})(\text{Backward Recursion}) \quad (4.14)$$

This is called Latent sequence given observations in HMM in the course notes

<sup>9</sup>A very good intuitive example from KPM is that we have a detective investigating a crime scene, and as we gather more observation, the prediction that we make tends to be more and more certain and less 'stochastic'.

### 4.3.1 Forward Filtering Recursion

$$\alpha_t(z_t) = p(z_t, x_{1:t}) = \sum_{z_{t-1}} p(z_{t-1}, z_t, x_{1:t}) \quad (4.15)$$

$$= \sum_{z_{t-1}} p(x_t | z_{t-1}, z_t, x_{1:t-1}) p(z_t | z_{t-1}, x_{1:t-1}) \underbrace{p(z_{t-1}, x_{1:t-1})}_{=\alpha_{t-1}(z_{t-1})} \quad (4.16)$$

$$\implies \alpha_t(z_t) = p(x_t | z_t) \sum_{z_{t-1}} p(z_t | z_{t-1}) \alpha_{t-1}(z_{t-1}) \quad (4.17)$$

Notice that our forward recursion contains our emission,  $p(x_t | z_t)$  and transition  $p(z_t | z_{t-1})$ . The base case of the recursion could be unwinded into

$$\alpha_1(z_1) = p(z_1, x_1) = p(z_1) p(x_1 | z_1) \quad (4.18)$$

### 4.3.2 Backward Filtering Recursion

$$p(x_{t+1:T} | z_t) = \sum_{z_{t+1}} p(z_{t+1}, x_{t+1:T} | z_t) \quad (4.19)$$

$$= \sum_{z_{t+1}} p(x_{t+2:T} | z_{t+1}, z_t, x_{t+1}) p(x_{t+1} | z_{t+1}, z_t) p(z_{t+1} | z_t) \quad (4.20)$$

$$\implies \beta_t(z_t) = \sum_{z_{t+1}} \underbrace{p(x_{t+2:T} | z_{t+1})}_{=\beta_{t+1}(z_{t+1})} p(x_{t+1} | z_{t+1}) p(z_{t+1} | z_t) \quad (4.21)$$

and hence if we recurse the above relationship, we will unwind to the base case

$$\beta_1(z_1) = p(x_{3:T} | z_2) p(x_2 | z_2) p(z_2 | z_1) \quad (4.22)$$

## 4.4 Sampling

The goal is to,

- Generate sample  $\{x^{(r)}\}_{r=1}^R$  from a probability distribution  $p(x)$ , and/or
- Estimate expectations of functions  $f(x)$  under some distribution  $p(x)$ , usually we want to estimate moments for a distribution

The word sample here refers to a single realization from a distribution rather than a set.

$$E = \mathbb{E}_{x \sim p(x)} [f(x)] = \int f(x) p(x) dx \quad (4.23)$$

### 4.4.1 Ancestral Sampling

**Generating marginal samples** If you are only interested in sampling a particular set of nodes, you can simply sample from all the nodes jointly, then ignore the nodes you don't need.

**Generating conditional samples** If you want to sample conditional on a node with no parents, that's also easy - you can simply do ancestral sampling starting from the nodes you have.

(#! **Important:**) However, to sample from a DAG conditional on leaf nodes is hard in the same way that inference is hard in general. E.g. sampling the unknown key in a crypto-system given the cypher-text but not knowing the plaintext. Finding ways to do this approximately is what a lot of the rest of the course will be about.

#### 4.4.2 Simple Monte Carlo

**(Definition) - Simple MC** Given  $\{x^{(r)}\}_{r=1}^R \sim p(x)$ , we want to estimate the expectation  $\mathbb{E}_{x \sim p(x)}[f(x)]$ .

$$E = \mathbb{E}_{x \sim p(x)}[f(x)] \approx \frac{1}{R} \sum_{r=1}^R f(x^{(r)}) = \hat{E} \quad (4.24)$$

#### Unbiasedness of MC

$$\mathbb{E}_{x \sim p(\{x^{(i)}\}_{r=1}^R)}[\hat{E}] = \mathbb{E} \left[ \frac{1}{R} \sum_{r=1}^R f(x^{(r)}) \right] \quad (4.25)$$

$$= \frac{1}{R} \sum_{r=1}^R \mathbb{E} [f(x^{(r)})] \quad (4.26)$$

$$= \frac{1}{R} \sum_{r=1}^R \mathbb{E}_{x \sim p(x)} [f(x)] \quad (4.27)$$

$$= \frac{R}{R} \mathbb{E}_{x \sim p(x)} [f(x)] \quad (4.28)$$

$$= E \quad (4.29)$$

■

#### Variance of MC

$$\text{var}[\hat{E}] = \text{var} \left[ \frac{1}{R} \sum_{r=1}^R f(x^{(r)}) \right] \quad (4.30)$$

$$= \frac{1}{R^2} \text{var} \left[ \sum_{r=1}^R f(x^{(r)}) \right] \quad (4.31)$$

$$= \frac{1}{R^2} \sum_{r=1}^R \text{var} [f(x^{(r)})] \quad (4.32)$$

$$= \frac{1}{R} \text{var} [f(x)] \quad (4.33)$$

We notice that the accuracy of MC estimates only depends on the variance of  $f$ , not on the dimension of  $x$ . Also, as the number of samples,  $R$ , increases, the variance  $\hat{E}$  will decrease at a rate proportional to  $\frac{1}{R}$ . ■

## 5 Lecture 7 - Stochastic Variational Inference (SVI/ADVI)

### 5.1 Motivation - Approximating Posterior Inference

Notice that if we want to calculate the *exact* posterior distribution, for example,

$$p(z|x) = \frac{p(x|z)}{p(x)} = \frac{p(x, z)p(z)}{\int p(x, z)dz} \quad (5.1)$$

the bottom integral over all possible states soon becomes intractable, making the computation of the entire posterior distribution intractable. Hence, instead, we will approximate the posterior inference with *variational methods*. Generally speaking, it works as follows

1. We introduce a family of distributions, with variational parameters  $\phi$ , denoted as  $q_\phi$ .
2. We want to encode a “distance metric” between  $p(z|x)$  and  $q_\phi(z)$ .
3. Try to minimize the proposed “distance metric”.

There are a few things to note

- First, as we will soon see, we use a measure of closeness that is **not** a distance metric. Also, using naïve distance metric could fail spectacularly (arbitrarily bad examples can be found if we use MSE).
- This formulation turns the Bayesian Inference into an optimization problem. If enough parts of the model is differentiable and could be well-approximated with Monte Carlo, then we can use gradient based optimization methods to solve this problem scalably.

### 5.2 The Kullback-Leibler Divergence

Suppose  $q_\phi$  denotes a family of distributions with its own **variational parameters**  $\phi$ , and  $p$  is some distribution. Then, the Kullback-Leibler divergence is defined as

$$D_{\text{KL}}(q_\phi(z|x) || p(z|x)) = \int q_\phi(z|x) \log \frac{q_\phi(z|x)}{p(z|x)} dz \quad (5.2)$$

$$= \mathbb{E}_{z \sim q_\phi} \left[ \log \frac{q_\phi(z|x)}{p(z|x)} \right] \quad (5.3)$$

**Properties** of KL Divergence

1.  $D_{\text{KL}}(q_\phi || p) \geq 0$ , non-negativeness
2.  $D_{\text{KL}}(q_\phi || p) = 0 \iff q_\phi = p$
3.  $D_{\text{KL}}(q_\phi || p) \neq D_{\text{KL}}(p || q_\phi)$ , i.e. KL Divergence is **not** a metric distance.

### 5.2.1 Expectation Propagation

If we try to reverse the arguments, it leads to a different kind of variational inference that is called “expectation propagation”. Typically, this leads to an algorithm that is more computationally expensive.

## 5.3 The Evidence Lower BOund (ELBO)

### 5.3.1 Definition and Equivalence Property

We cannot minimize the KL divergence exactly, thus we introduce a function, ELBO, and try to minimize that. Notice that ELBO is equal to the actual KL divergence up to a constant.

$D_{\text{KL}}(q_\phi||p)$  is intractable because it contains  $p(z|x)$ , which is intractable.

$$D_{\text{KL}}(q_\phi(z|x)||p(z|x)) = \mathbb{E}_{z \sim q_\phi} \log \frac{q_\phi(z|x)}{p(z|x)} \quad (5.4)$$

$$= \mathbb{E}_{z \sim q_\phi} \left[ \log \left( q_\phi(z|x) \cdot \frac{p(x)}{p(z, x)} \right) \right] \quad (5.5)$$

$$= \mathbb{E}_{z \sim q_\phi} \log \frac{q_\phi(z|x)}{p(z, x)} + \mathbb{E}_{z \sim q_\phi} \log p(x) \quad (5.6)$$

$$= -\mathcal{L}(\phi; x) + \log p(x) \quad (5.7)$$

where  $\mathcal{L}(\phi; x) = -\mathbb{E}_{z \sim q_\phi} \log \frac{q_\phi(z|x)}{p(z, x)}$  is defined as the ELBO. If we re-arrange what we have above,

$$\log p(x) = \mathcal{L}(\phi; x) + D_{\text{KL}}(q_\phi(z|x)||p(z|x)) \quad (5.8)$$

and we notice that  $D_{\text{KL}}(q_\phi(z|x)||p(z|x)) \geq 0$ , it then follows that

$$\mathcal{L}(\phi; x) \leq \log p(x) \quad (5.9)$$

Hence, to minimize the KL divergence we need to maximize the ELBO, and they are equivalent.

### 5.3.2 Optimization goal of ELBO

We have, by definition of ELBO, that

$$\mathcal{L}(\phi; x) = -\mathbb{E}_{z \sim q_\phi} \log \frac{q_\phi(z|x)}{p(x, z)} \quad (5.10)$$

$$= \mathbb{E}_{z \sim q_\phi} [\log p(x, z) - \log q_\phi(z|x)] \quad (5.11)$$

Recall, that our goal is to use gradient based methods to optimize this objective, meaning that we have to somehow compute, estimate to be precise in this case, the gradient of ELBO  $\mathcal{L}(\phi; x)$

$$\nabla_\phi \mathcal{L}(\phi) = \nabla_\phi \mathbb{E}_{z \sim q_\phi(z|x)} [\log p(x, z) - \log q_\phi(z|x)] \quad (5.12)$$

and we know from a previous lecture that one can get an unbiased estimator of any expectation as long as we can sample from the distribution and evaluate the function, using simple Monte Carlo.



### 5.3.3 Pathwise Gradient

Above, we saw that the objective is the gradient over an expectation and we now would like to make it an expectation over gradient. In general, such swapping is not allowed and we can only do so if the distribution we're taking the expectation over **does not** depend on the parameter, i.e.  $\phi$ . Our goal is now to factor out the randomness of from  $q$ , and put it into a parameterless, fixed source of noise  $p(\varepsilon)$ . Formally, we need to find function  $T(\phi, \varepsilon)$  such that

$$\left. \begin{array}{l} \varepsilon \sim p(\varepsilon) \\ z = T(\phi, \varepsilon) \end{array} \right\} \implies z \sim q_\phi(z) \quad (5.13)$$

We usually start with  $p(\varepsilon)$  being uniform or normal. Indeed it is not always easy to find these functions<sup>10</sup> and as a easy concrete example, we know

$$\left. \begin{array}{l} \varepsilon \sim \mathcal{N}(\varepsilon|0, 1) \\ z = \sigma\varepsilon + \mu \end{array} \right\} \implies z \sim \mathcal{N}(z|\mu, \sigma) \quad (5.14)$$

Then, by applying the trick we know,

$$\nabla_\phi \mathcal{L}(\phi) = \nabla_\phi \mathbb{E}_{z \sim q_\phi(z|x)} [\log p(x, z) - \log q_\phi(z|x)] \quad (5.15)$$

$$= \nabla_\phi \mathbb{E}_{\varepsilon \sim p(\varepsilon)} [\log p(x, T(\phi, \varepsilon)) - \log q_\phi(T(\phi, \varepsilon)|x)] \quad (5.16)$$

$$= \mathbb{E}_{\varepsilon \sim p(\varepsilon)} \nabla_\phi [\log p(x, T(\phi, \varepsilon)) - \log q_\phi(T(\phi, \varepsilon)|x)] \quad (5.17)$$

## 5.4 Tutorial

### 5.4.1 Derivation Using Jensen's Inequality

Recall that by when Jensen's inequality is applied to probability distributions, we know that when  $f$  is concave,

$$f(\mathbb{E}[X]) \geq \mathbb{E}[f(X)] \quad (5.18)$$

Then, for the log probability of the observations,

$$\log p(x) = \log \int p(x, z) dz \quad (5.19)$$

$$= \log \int p(x, z) \frac{q_\phi(z|x)}{q_\phi(z|x)} dz \quad (5.20)$$

$$= \log \mathbb{E}_{z \sim q_\phi} \frac{p(x, z)}{q_\phi(z|x)} \quad (5.21)$$

where we apply Jensen's Inequality

$$\implies \mathbb{E}_{z \sim q_\phi} \frac{p(x, z)}{q_\phi(z|x)} \geq \mathbb{E}_{z \sim q_\phi} \log \frac{p(x, z)}{q_\phi(z|x)} \quad (5.22)$$

$$= - \mathbb{E}_{z \sim q_\phi} \log \frac{q_\phi(z|x)}{p(x, z)} \quad (5.23)$$

$$= \mathcal{L}(\phi; x) \quad (5.24)$$

---

<sup>10</sup><http://blog.shakirm.com/2015/10/machine-learning-trick-of-the-day-4-reparameterisation-tricks/>

### 5.4.2 Interpretations of ELBO

For the ELBO defined as

$$\text{ELBO} \triangleq \mathcal{L}(\phi; x) = - \mathbb{E}_{z \sim q_\theta} \log \frac{q_\phi(z|x)}{p(x, z)} \quad (5.25)$$

we have

1. The most general interpretation is just breaking up the logarithm. We have

$$\mathcal{L}(\phi; x) = - \mathbb{E}_{z \sim q_\phi} \log \frac{q_\phi(z|x)}{p(x, z)} \quad (5.26)$$

$$= \mathbb{E}_{z \sim q_\phi} \log \frac{p(x, z)}{q_\phi(z|x)} \quad (5.27)$$

$$= \mathbb{E}_{z \sim q_\phi} \log \frac{p(z)p(x|z)}{q_\phi(z|x)} \quad (5.28)$$

$$= \mathbb{E}_{z \sim q_\phi} [\log p(x|z) + \log p(z) - \log q_\phi(z|x)] \quad (5.29)$$

2. We can try to write this using entropy as

$$\mathcal{L}(\phi; x) = \mathbb{E}_{z \sim q_\phi} [\log p(x|z) + \log p(z)] \mathcal{H}(q_\phi(z|x)) \quad (5.30)$$

3. We can frame ELBO as a trade off

$$\mathcal{L}(\phi; x) = \underbrace{\mathbb{E}_{z \sim q_\phi} [\log p(x|z)]}_{(\dagger)} - \underbrace{D_{\text{KL}}(q_\phi(z|x) \| p(z))}_{(\ddagger)} \quad (5.31)$$

- $(\dagger)$  is “reconstruction likelihood”, i.e. how probable is  $x$  given  $z$ , which encourages the model to choose the distribution which best reconstructs the data.
- $(\ddagger)$  acts as regularization, enforcing the idea that our parametrization shouldn't move us too far from the true distribution.

## 6 Lecture 8 - Sampling and Monte Carlo Methods

### 6.1 Motivation: Sampling

Note: In this section, we refer to a “sample” as a single realization  $x$  that is from a probability distribution  $p(x)$ .<sup>11</sup>

**Problems to be solved** Monte Carlo methods are computational techniques that make use of random numbers, it aims to solve one or both of the following problems

1. To generate samples

$$\{x^{(r)}\}_{r=1}^R \sim p(x) \quad (6.1)$$

2. and to estimate expectations of functions  $(\phi(x))$  under the distribution of  $p(x)$ , i.e.,

$$\Phi = \mathbb{E}_{x \sim p(x)} [\phi(x)] = \int \phi(x) p(x) dx \quad (6.2)$$

---

<sup>11</sup>This is different from the normal sense of sample, which means one batch of sample from a distribution.

**Sampling  $p(x)$  is like hornets' nest** To sample from

$$p(x) = \frac{\tilde{p}(x)}{Z} \quad (6.3)$$

there are two main complications even if we know how to evaluate  $\tilde{p}(x)$

1. We don't typically know the normalizing constant,  $Z$ .
2. Even if we did know  $Z$ , the problem of drawing samples from  $p(x)$  is still a challenging one, especially in high-dimensional spaces, because there is no obvious way to sample from  $p$  without enumerating most or all of the possible states.

Should this be  $\tilde{p}$  ??

Note: There are only a handful of distribution, for example gaussian distribution, that we can easily sample from in a high dimensional space.

## 6.2 Lattice Discretization

Imagine that we wish to draw samples from the density  $p(x) = \frac{\tilde{p}(x)}{Z}$ , then what we can do is discretize the variable  $x$  and sample from the discrete probability distribution over a finite set of uniformly spaced points  $S = \{x_i\}_i$ . Then for  $x_i \in S$ , we evaluate  $\tilde{p}(x_i)$  after which we can compute the normalizer

$$Z = \sum_i \tilde{p}_i(x_i) \quad (6.4)$$

and

$$p_i = \frac{\tilde{p}_i}{Z} \quad (6.5)$$

The main problem with this method, however, comes from the accuracy and cost of evaluating at lattice. If we want to sample  $D$  uniformly spaced points in one dimension and the system has  $N$  dimensions in total. Then, the the complexity is asymptotically  $\mathcal{O}(D^N)$  which grows exponentially as dimension grows and soon become intractable.

## 6.3 Importance Sampling

Note: Importance Sampling doesn't help us solve the problem of generating samples from  $p(x)$ , rather it is a method that helps us estimate expectation of a function  $\phi(x)$ , i.e. problem 2.

### Assumptions

1. To do importance sampling, we still need to assume that we can evaluate  $p(x)$  (density that we wish to draw samples from) within a multiplicative constant. To be precise, we can evaluate a function  $\tilde{p}(x)$  such that

$$p(x) = \tilde{p}(x)/Z \quad (6.6)$$

2. We further assume that we have a simpler density,  $q(x)$  from which it is easy to **sample** from, i.e.,  $x \sim q(x)$  and easy to **evaluate**, i.e.,  $\tilde{q}(x)$

$$q(x) = \frac{\tilde{q}(x)}{Z_q} \quad (6.7)$$

where the density  $q(x)$  is referred to as the sampler density.

**Sampling Procedure** In importance sampling, we generate in total  $R$  samples from  $q(x)$ , i.e.  $\{x^{(r)}\}_{r=1}^R \sim q(x)$ . If these points were samples from  $p(x)$  then we could estimate  $\Phi$  by

$$\Phi = \mathbb{E}_{x \sim p(x)}[\phi(x)] \approx \frac{1}{R} \sum_{r=1}^R \phi(x^{(r)}) = \hat{\Phi} \quad (6.8)$$

i.e., we could use a *Simple Monte Carlo Estimator*. The problem is that these sample are from  $q$  rather than  $p$  which can be completely unrelated densities. To solve this over-represented/under-represented problem, we can take into account that we are sampling from a different distribution by introducing weights, (essentially corrections)

$$\tilde{w}_r = \frac{\tilde{p}(x^{(r)})}{\tilde{q}(x^{(r)})} \quad (6.9)$$

Finally, we re-write our estimator under  $q$ , i.e.

$$\Phi = \int \phi(x) p(x) dx \quad (6.10)$$

$$= \int \phi(x) \frac{p(x)}{q(x)} q(x) dx \quad (6.11)$$

$$\approx \frac{1}{R} \sum_{r=1}^R \phi(x^{(r)}) \frac{p(x^{(r)})}{q(x^{(r)})} \quad (6.12)$$

$$= \frac{Z_q}{Z_p} \frac{1}{R} \sum_{r=1}^R \phi(x^{(r)}) \cdot \frac{\tilde{p}(x^{(r)})}{\tilde{q}(x^{(r)})} \quad (6.13)$$

$$= \frac{Z_q}{Z_p} \frac{1}{R} \sum_{r=1}^R \phi(x^{(r)}) \cdot \tilde{w}_r \quad \text{since } w_r = \frac{\tilde{w}_r}{\sum_{r=1}^R \tilde{w}_r} \quad (6.14)$$

$$= \frac{\frac{1}{R} \sum_{r=1}^R \phi(x^{(r)}) \cdot \tilde{w}_r}{\frac{1}{R} \sum_{r=1}^R \tilde{w}_r} \quad \text{since } \frac{Z_p}{Z_q} = \frac{1}{R} \sum_{r=1}^R \tilde{w}_r \quad (6.15)$$

$$= \frac{1}{R} \sum_{r=1}^R \phi(x^{(r)}) \cdot w_r \quad (6.16)$$

$$= \hat{\Phi}_{iw} \quad (6.17)$$

The estimator  $\hat{\Phi}_{iw}$  is biased, but is consistent.

**Pitfalls** As it is easy to see, if we choose  $q$  to be some short tailed distribution, we might not get a good coverage for samples on  $p$ , and thus making the estimation not accurate. Notice that we want to support the entire real line with this method<sup>12</sup> and hence we cannot use a naïve uniform over a certain range, and using uniform over the entire real line will cause uniform distribution pdf to be infinitely small everywhere.

---

<sup>12</sup>Or otherwise this will just be uniform sampling. Indeed, importance sampling is a generalization of uniform sampling.

**Remedy** As a way to (semi-) solve the above pitfall, we can choose  $q$  to be a heavy-tailed distribution, for example Cauchy. Of course, the best of all solution is to find a  $q$  that looks like  $p$ , but this is in a lot of cases an un achievable goal.

## 6.4 Rejection Sampling

In rejection sampling we assume again a one-dimensional density  $p(x) = \tilde{p}(x)/Z$  that is too complicated a function for us to be able to sample from it directly. We assume that we have a simpler **proposal** density  $q(x)$  which we can evaluate (within a multiplicative factor  $Z_q$ , as before), and from which we can generate samples. We further assume that we know the value of a constant  $c$  such that

$$c\tilde{q}(x) > \tilde{p}(x) \quad \text{for some } c \text{ constant}, \forall x \quad (6.18)$$

It is somewhat important to note that in rejection sampling we care only about the proportion, and we need a  $q$ , or  $\tilde{q}$  actually, that “covers”  $p$ , or  $\tilde{p}$  actually. This means that we don’t care, and thus don’t need to know the normalizing constant to perform the sampling procedure. Naïvely, we can just scale  $q$  to a huge  $\tilde{q}$ , and as long as it covers  $\tilde{p}$ , we are in a good shape.

**Sampling Procedure** The procedure works as follows

1. Generate two random numbers, the first  $x$  is such that  $x \sim q(x)$ , and the second  $u$  is generated such that  $u \sim \text{Unif}[0, c\tilde{q}(x)]$
2. Evaluate  $\tilde{p}(x)$  and accept or reject the sample  $x$  by comparing the value of  $u$  with the value of  $\tilde{p}(x)$ , to be more specific
  - If  $u > \tilde{p}(x)$ , then we reject  $x$ ,
  - otherwise we accept  $x$ , effectively adding  $x$  to our set of samples. (value of  $u$  is discarded, and in the next iteration  $u$  will be sampled again from, possibly, a different uniform distribution. )

### Pitfalls

- Rejection sampling works the best if we are able to find  $q$  that is a good approximation to  $p$ . If  $q$  is very different from  $p$  then, for  $cq$  to exceed  $p$  everywhere,  $c$  will almost always necessarily be large. This will cause us a very high rejection ratio, and making the algorithm super inefficient.
- In high-dimensional space, the problem is even worse and it is very likely that the requirement that  $c\tilde{q} > \tilde{p}$  will force  $c$  to be so huge that acceptances will be very rare indeed.
- Besides, finding such a value of  $c$  may be difficult too, since in many problems we know neither where the modes of  $\tilde{p}$  are located nor how high they are.
- In general,  $c$  grows exponentially with the dimensionality  $N$ , so the acceptance rate is expected to be exponentially small in  $N$ , since

$$\text{acceptance rate} = \frac{\text{area under } \tilde{p}}{\text{area under } c\tilde{q}} = \frac{1}{Z} \quad (6.19)$$

was ist das?

## 6.5 Metropolis-Hasting Method: Example of MCMC

The proposed importance sampling and rejection sampling work well only if the proposal density  $q(x)$  is similar to  $p(x)$ . In general, in a high dimensional case, such  $q$  is very hard to find. In Metropolis-Hasting, in contrast to importance and rejection sampling methods, it is not necessary for  $q(x'|x^{(t)})$  to look similar to  $p(x)$  in order for the algorithm to be useful. The density  $q(x'|x^{(t)})$  might be a simple distribution such as a gaussian centred at  $x^{(t)}$ , but can in general be any *fixed* density from which we can draw samples.

### Sampling Procedure

- A tentative new state  $x'$  is generated from the proposal density  $q(x'|x^{(t)})$ , then we compute

$$a = \frac{\tilde{p}(x') q(x^{(t)}|x')}{\tilde{p}(x^{(t)}) q(x'|x^{(t)})} \quad (6.20)$$

- If  $a \geq 1$ , then the new state is accepted.
- Otherwise, the new state is accepted with probability  $a$ . Then, we perform update as follows,

$$x^{(t+1)} \leftarrow \text{accepted ? } x' : x^{(t)} \quad (6.21)$$

**Note on Difference** Notice that in rejection sampling, rejected points are discarded and have no influence on the list of samples  $x_r^{(r)}$  that we collected. Here, a rejection causes the current state to be written again onto the list.

**Convergence and MCMC** Metropolis-Hastings converges to  $p(x)$  for any  $q(x'|x^{(t)}) \geq 0 \quad \forall x', x^{(t)}$  as  $t \rightarrow \infty$ . That is, our list of samples is such that

$$\{x^{(r)}\}_{r=1}^R \rightarrow p(x) \quad (6.22)$$

The Metropolis-Hastings method is an example of a Markov Chain Monte Carlo method. In contrast to rejection sampling, where the accepted points are independent samples from the desired distribution, MCMC methods involve a Markov process in which a sequence of states is generated, each sample  $x^{(t)}$  having a probability distribution that depends on the previous value,  $x^{(t-1)}$ . Since successive samples are dependent, the Markov chain may have to run for a considerable time in order to generate samples that are effectively independent samples from  $p$ .

## 7 Lecture 9 - Amortized Inference and Variational Auto Encoders

### 7.1 Motivations

#### 7.1.1 Efficient Use of (Accumulative) Partial Evidence

With enough experience, doctors, plumbers, detectives, etc. can very quickly tell what is going on and what they still need more information about, if they've seen enough similar

situations. Hence, perhaps we could somehow train a neural network to look at the data for a person  $x_i$ , and then output an approximate posterior  $q_{\phi_i}(z_i|x_i)$

## 7.2 Amortized Inference

Here “amortized” essentially means to spread out over time. Previously, we do SVI from scratch every time we see a new datapoint. Now, instead, we will gradually learn a function that can look at the data for a person  $x_i$ , and then output an approximate posterior  $q_{\phi}(z_i|x_i)$ . This is known as a **recognition model**, instead of a separate  $\phi_i$  for each data example, we will just have a single global  $\phi$  that specifies the parameter of the recognition model. Because the relationship between data and posteriors is complex and hard to specify by hand, we will do this through a neural net. (we make the neural net learn the function!)

### 7.2.1 Simple Example (Parametrizing Gaussian)

We have seen one way to specify a probability distribution given an input with neural networks in assignment 2: we can simply have a network take in  $x_i$ , and output the mean and variance vector for a Gaussian, parametrized as

$$q_{\phi}(z_i|x_i) = \mathcal{N}(z_i|\mu_{\phi}(x_i), \Sigma_{\phi}(x_i)) \quad (7.1)$$

The graphical model for such model specification has  $\phi$  **outside of** the plate, which means we only have one global  $\phi$  to update and optimize.

### 7.2.2 Amortized Inference - Algorithmic Procedure

1. Sample a datapoint
2. compute parameters of approximate posterior (recognition)
3. compute gradient of Monte Carlo estimate of ELBO with respect to  $\phi$ , i.e.  $\nabla_{\phi}^{\sim MC} ELBO$
4. update according to gradient descent

Then, when we want to make predictions about a new datapoint, we can use our fast recognition model if we want. Of course, we might also want to stop and do something slower but more accurate, for example like per-sample SVI, and/or MCMC.

### 7.2.3 Optimizing Model Parameters

Now that  $p_{\theta}(x)$  depends on parameters  $\theta$ , then the ELBO is a function of both  $\phi$  and  $\theta$ . We can optimize them together, still, using the re-parametrize trick we saw earlier.

$$\nabla_{\theta, \phi} \mathcal{L}(\phi) = \nabla_{\theta, \phi} \mathbb{E}_{z \sim q_{\phi}(z|x)} [\log p_{\theta}(x, z) - \log q_{\phi}(z|x)] \quad (7.2)$$

$$= \mathbb{E}_{\epsilon \sim p(\epsilon)} \nabla_{\theta, \phi} [\log p_{\theta}(x, T(\phi, \epsilon)) - \log q_{\phi}(T(\phi, \epsilon)|x)] \quad (7.3)$$

This allows us to jointly fit the model parameters and the recognition network, by subsampling training examples and using simple Monte Carlo with gradient descent optimizers. Such formulation is referred to as a variational auto encoder, VAE (See Figure 7.1 for layout of architecture).

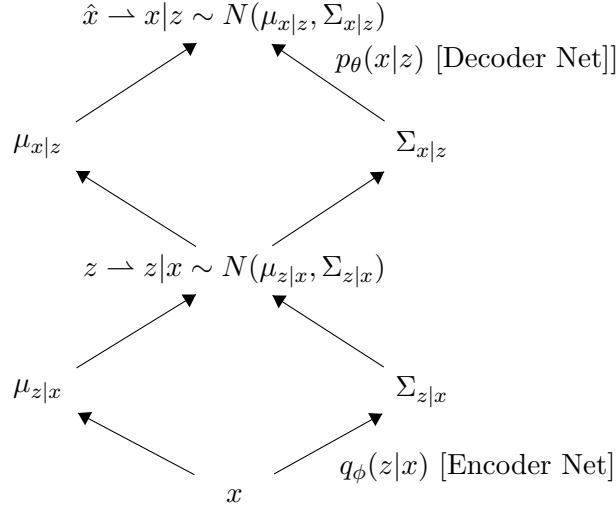


Figure 7.1: Variational Autoencoder Architecture

### 7.3 Variational Auto Encoders (VAEs)

#### 7.3.1 Autoencoders: Definition and Problems

An autoencoder takes an input  $x$ , and encodes it into a code vector  $z$ <sup>13</sup>, after which the code vector  $z$  is decoded into  $\hat{x}$  that is “similar” to  $x$ , our original data. Essentially, what we want to learn is two mappings, namely the encoder  $g(x) \rightarrow z$  and the decoder  $f(z) \rightarrow \hat{x}$ . In general,  $g$  and  $f$  could be arbitrarily complex, and thus we model them neural nets. The optimization goal of such problem is to minimize the reconstruction error, i.e.

$$\mathcal{J} = \text{reconstruction\_error}(x, \hat{x}) \quad (7.4)$$

**Unsupervised Nature of Autoencoders** Notice that here the cost defined in equation 7.4 doesn’t rely on *any* explicit label input. Rather, it is a form of unsupervised learning, where the learning goal is to learn the mappings to reconstruct the inputs, i.e.  $\tilde{x} = f \circ g(x) \cong x$ .

#### 7.3.2 Problems with Deterministic Autoencoders

**Proximity in data space doesn’t mean proximity in feature space** The codes that we learn by the model is deterministic, i.e.

$$\begin{aligned} g(x_1) = z_1 &\Rightarrow f(z_1) = \tilde{x}_1 \\ g(x_2) = z_2 &\Rightarrow f(z_2) = \tilde{x}_2 \end{aligned} \quad (7.5)$$

but proximity in feature space is not enforced for inputs in close proximity in data space, i.e.

$$z_1 \approx z_2 \not\Rightarrow x_1 \approx x_2 \quad (7.6)$$

<sup>13</sup>Notice that for autoencoders to actually “encode” the inputs, we necessarily have to choose the dimension of  $z$  so that it forms a bottleneck. If we allow  $z$  to have the same (or even larger) dimensionality with/than  $x$  and  $\hat{x}$ , then likely the model will learn two identity mappings!



If the space has regions where no data gets encoded to, and you sample/generate a variation from there, the decoder will generate an unrealistic output, because the decoder has no idea how to deal with that region of the latent space. During training, it never saw encoded vectors coming from that region of latent space

**Adding Noise By Hand is Hard** We said that we wish to embed  $x$ 's into a lower dimension space, and code it as  $z$ . However, in practice it is very hard to determine what the low dimension  $z$  should be. Similar to this,

1. We can add noise to data before encoding, after which we try to reconstruct the original data. But how much noise?
2. We can add noise to the latent variable, code,  $z$  after encoding, after which we try to reconstruct the original data. Again, how much noise should we add?

### 7.3.3 Variational Auto Encoders

In an hope to address the problems above, we introduce stochasticity into our model: even for the same input, while the mean and standard deviations remain the same, the actual encoding will somewhat vary on every single pass simply due to sampling. Figure 7.1 illustrates the architecture of a variational autoencoder.

**VAE Solves Problems in Section 7.3.2, WHY?** The VAE generation model learns to reconstruct its inputs not only from the encoded points but also from the area around them. This allows the generation model to generate new data by sampling from an “area/volume” instead of only being able to generate already seen data corresponding to the particular fixed encoded points.

### 7.3.4 Utilization of Latent Variable Models

The question is: once we have obtained the model, what do we do with it?

If a latent variable model has a compact prior and vector valued code  $z$ 's, we can

1. Sample new data to check the model,
2. Encode (sample from the approximate posterior) of two data points, and interpolate between them in the latent space, then decode along the path
3. Learn a function to predict some property of the examples from the low-dimensional  $z$  space instead of directly from  $x$ . (semi-supervised learning, essentially use unsupervised learning to cleanup the data, and to learn a “better” representation.)