

Machine Learning Summer School 2018

Miquel Perello-Nieto

August 27–September 7, 2018

Contents

Contents	iii
Acronyms	vii
1 Introduction	1
2 Bayesian Deep Learning: Planting the seeds of probabilistic thinking by Shakir Mohamed	3
2.1 Introduction 27 Aug. Mon. 9:30–11:00	3
2.1.1 Learning and inference	4
2.2 Inferential questions: Mon. 11:30–13:00	5
2.2.1 Hutchinson’s Trick	6
2.2.2 Probability Flow Tricks	6
2.2.3 Stochastic Optimization	6
2.2.4 Pathwise Estimator	6
2.2.5 Log-derivative trick	7
2.2.6 Score-function estimator	7
2.2.7 Evidence Bounds	7
2.2.8 Density Ratio Trick	7
2.3 Unsupervised Learning Thu. 9:30–11:00	7
2.3.1 Fully-observed models	8
2.3.2 Latent variable models	8
2.3.3 Implicit Models	9
2.4 Model-Inference-Algorithm	9
2.5 Variational Inference	9
2.6 Mean-Fields	9
2.7 Variational Optimisation	10
2.8 Reinforcement learning as generative models	10
3 An Introduction to Gaussian Processes by Richard Turner	11
3.1 Covariance functions	12
3.1.1 References	13
3.2 Using Gaussian Processes: Models, applications and connections 16:30–18:00	13
3.2.1 Deep Gaussian Processes	15
3.2.2 References	15
3.3 Large data and non-linear models Tue. 11:30–13:00	15
3.4 A brief history of gaussian process approximation	16
3.4.1 Fully independent training conditional (FITC) approximation	16
3.4.2 Variational free-energy method (VFE)	17
4 Deep Learning by Rob Fergus	19
4.1 Deep Supervised Learning	19
4.1.1 History of Neural Nets	19
4.1.2 Deep learning vs traditional approaches	19
4.1.3 Some issues with Deep Learning	20
4.1.4 Convolutional Neural Networks	20

4.1.5	Stochastic Gradient Descent	20
4.1.6	Some practical debugging tips from M. Ranzato	20
4.1.7	Weights' initialization	20
4.1.8	Deep Residual Learning	21
4.1.9	Weakly supervised pretraining	21
4.2	Unsupervised Learning	21
4.2.1	Self supervised learning	21
4.2.2	Auto-Encoder	21
4.2.3	Variational Auto-Encoder	21
4.2.4	Generative Adversarial Networks	21
4.2.5	Stacked Auto-Encoders	22
4.2.6	Other approaches	22
4.2.7	Application examples	22
4.3	Deep Learning Models Rob Fergus 9:30–11:00	22
4.3.1	Memory in Deep networks	22
4.4	Deep Nets for sets	23
	Bibliography	23
5	Statistical machine learning and convex optimization by Francis Bach	25
5.1	Introduction 14:30–18:00	25
5.2	Classical methods for convex optimization	25
5.2.1	Lipschitz continuity	26
5.2.2	Smoothness and strong convexity	26
5.2.3	Analysis of empirical risk minimization	26
5.2.4	Accelerated gradient methods (Nesterov, 1983)	26
5.2.5	Optimization fro sparsity-inducing norms	26
5.2.6	Summary about minimization of convex functions	26
5.3	Convex stochastic approximation	27
5.4	Summary of rates of convergence	27
5.5	Conclusions	27
	Bibliography	28
6	Supervised Learning and Text Classification by Kyunghyun Cho	29
6.1	Introduction to supervised learning with ANNs Wed. 11:30–13:00	29
6.1.1	Loss minimization	29
6.2	Text classification	30
6.2.1	How to represent a sentence	30
6.3	Natural Language Models	31
6.3.1	Autoregressive language modelling	31
6.3.2	N-Gram language models	31
6.3.3	Neural N-Gram Language Model	31
6.3.4	Convolutional Language Models	32
6.3.5	CBoW Language Models (infinite context)	32
6.3.6	Recurrent Language Models	32
6.3.7	Recurrent Memory Networks	32
6.4	Recurrent Networks and Backpropagation	32
6.4.1	Gated recurrent units	32
6.4.2	Lessons from GRU/LSTM	32
6.5	Neural Machine Translation 14:30–16:00	33
6.5.1	History of machine translation	33
6.5.2	Encoding: Token representation	33
6.5.3	Decoding: conditional language modeling	33
6.5.4	In practice	33
6.6	Current and ongoing projects	33
6.6.1	Multilingual translation	34
6.6.2	Real-Time Translation (learning to decode)	34
	Bibliography	34

7 Causality by Joris Mooij	37
7.1 Introduction	37
7.2 Defining causality in terms of probabilities	37
7.3 Causal Inference: Predicting Causal Effects	38
7.4 Resolving Simpson's paradox	38
7.5 Causal Discovery: from data to causal graph	39
7.5.1 Local Causal Discovery (LCD)	39
7.6 Practical application	39
7.7 Conclusions	39
Bibliography	40
8 Reinforcement Learning by Jan Peters	41
8.1 Optimal Control Systems	41
8.1.1 Markov Decision Problems	41
8.1.2 Basic reinforcement learning loop	41
8.1.3 Linear Quadratic Gaussian Systems	42
8.2 Value Function Methods	42
8.2.1 Markov Decision Processes (MDP)	43
8.2.2 Temporal difference learning	43
8.2.3 Approximating the value function	43
8.2.4 Batch-Mode Reinforcement Learning	43
8.3 Policy Search	44
8.3.1 Black-box approaches	44
8.3.2 Likelihood-Ratio Policy Gradient methods	44
8.4 Key problems	44
Bibliography	44
9 Probabilistic Numerics: Nano-machine-learning by Michael A Osborne	45
9.1 Upper confidence bound	47
9.2 Information-theoretic methods	48
9.3 Technology at work: The future of automation Tue. 9:30–11:00	48
Bibliography	48
10 Kernel methods by Arthur Gretton	49
10.1 Reproducing Hilbert spaces	49
10.2 Interlude: divergence measures	49
10.3 Two-sample testing with MMD	50
10.4 Training GANs with MMD Wed. 9:30–11:00	50
10.4.1 The kernel inception distance (KID)	50
10.5 Testing statistical dependence	50
10.5.1 Finding covariance with smooth transformations	50
10.5.2 Application: dependence detection across languages	50
Bibliography	51
11 An introduction to Bayesian nonparametrics by Sinead Williamson	53
11.1 The Dirichlet process	53
11.1.1 An urn representation	53
11.1.2 Exchangeability	53
11.1.3 Choosing the number of clusters	53
11.1.4 Constructing an appropriate prior	53
11.2 Dirichlet process and Dirichlet marginals	54
11.2.1 Conjugacy of the multinomial	54
11.2.2 The Chinese restaurant process	54
11.2.3 The stick breaking construction	54
11.2.4 Indian Buffet Process	54
11.2.5 Building latent feature models using the IBP	55
11.2.6 Summary	55

11.2.7 Latent Dirichlet allocation	55
11.2.8 Hierarchical Dirichlet process	55
11.2.9 The Chinese restaurant franchise	56
11.2.10 Basic network models: Erdős-Renyi models	56
11.3 Further resources	56
Bibliography	56
12 Machine Learning and Causal Inference for (Reliable) Decision Support. by Suchi Saria	57
12.1 Observed confounders	58
12.1.1 Feature Matching	58
12.2 References	58
12.3 Unstable paths	59
Bibliography	59
13 Machine Learning in the industry	61
13.1 Real-world ML Challenges at the Scale of Banking by BBVA Data and Analytics	61
13.2 Data Efficient Reinforcement Learning by PROWLER.io	61
13.3 Lynx: real-time accurate fraud detection over massive data. Instituto de Ingenieria del conocimiento by Álvaro Barbero Jiménez	61
13.4 Microsoft Research by Sebastian Nowozin	62
13.4.1 Timelines	62
Bibliography	63
14 Generative Adversarial Networks by Sebastian Nowozin	65
14.1 Probabilistic models	65
14.2 Example applications of GANs	65
14.3 Principles of estimation	66
14.3.1 Implicit models	66
14.4 GAN models	66
14.4.1 Divergences and f-GAN family	66
14.4.2 IPM	67
14.4.3 IPM family: MMD	67
14.4.4 IPM family: Wasserstein GANs	68
14.5 Problems and Fixes: Mode Collapse, Instability	68
14.5.1 Spectral Normalization	68
14.6 Implicit models more generally	69
14.7 Open research problems	69
Bibliography	69
15 Advances in Machine Learning for Molecules by José Miguel Hernández-Lobato	71
15.1 Molecular representation for ML	72
15.1.1 Molecular fingerprints	72
15.1.2 SMILES	73
15.1.3 Graph Neural Networks (GNNs)	74
Bibliography	75
Glossary	77

Acronyms

ANN Artificial Neural Network. 46, 48, 50, 51

GNN Graph Neural Network. 89–91

NLP Natural Language Processing. 90

WLN Weisfeiler-Lehman Network. 91

Chapter 1

Introduction

- 150 out of 500
- Aug 29 19:30 guided tour through downtown Madrid
- Sep 1 9:00am Visti to segovia
- Sep 5 20:00

Chapter 2

Bayesian Deep Learning: Planting the seeds of probabilistic thinking by Shakir Mohamed

- Shakir Mohamed - Cambridge, CIFAR, DeepMind

2.1 Introduction 27 Aug. Mon. 9:30–11:00

- Different views of probability
 - Statistical probability: frequency ratio of items
 - Logical Probability: Degree of confirmation of an hypothesis based on logical analysis
 - Probability as Propensity: probability used for predictions
 - Subjective probability: probability as a degree of belief
 - * Probability is a measure of a belief in a proposition given evidence. A description of a state of knowledge.
 - * Different observers with different information will have different beliefs
- $p(x)$ probability of x , $p^*(x)$ true probability of x
- Conditions: $p(x) \geq 0$, $\int p(x)dx = 1$
- Bayes rule: $p(z|x) = \frac{p(x|z)p(z)}{p(x)}$
- Parametrisation: $p_\theta(x|z) = p(x|z; \theta)$
- Gradient: missed
- Generalised Linear Regression

$$\mu = w^T x + b \tag{2.1}$$

- Recursive Generalised Linear Regression
- Recursively compose the basic linear functions
- Gives a deep neural network
- Probabilistic model
- $p(y|x) = p(y|h(x); \theta)$

- likelihood function (log of the previous one)
- $\mathcal{L}(\theta) = \sum_n \log p(y_n|x_n; \theta)$
- Sometimes the likelihood is not computationally tractable because of an integral
- Likelihoods can give you estimates of parameters
- They are efficient estimators, are asymptotically unbiased
- It is possible to make statistical tests on the likelihood, this can give you confidence intervals
- Pool information: combine the outcomes of different data sources
- Biggest problem: misspecification, inefficient estimates or confidence intervals/test can fail completely (assuming Gaussian while data is not)

A likelihood function with regularization term

$$\mathcal{L}(\theta) = \sum_n \log p(y_n|x_n; \theta) + \frac{1}{\lambda} \mathcal{R}(\theta) \quad (2.2)$$

Other names for the regularization are ...

The Maximum a Posteriori estimation (MAP)

Shows an example of two instantiations for a MAP estimate where the answer changes from 0 to 1. This problem arises because of the variable scale. In order to solve that, we need to learn more than just the mean.

In Bayesian analysis there are two core components: the evidence

$$p(y|x) = \int p(y|h(x); \theta) p(\theta) d\theta \quad (2.3)$$

and the posterior

$$p(\theta|y, x) \propto p(y|h(x); \theta) p(\theta) \quad (2.4)$$

In Bayesian analysis, all the things that are not observed need to be integrated over (averaged out)

Intractable Integrals do not have a closed form, or very high-dimensional quantities that can't be computed (e.g. using quadrature)

2.1.1 Learning and inference

In statistics there is no distinction between learning and inference, only inference (or estimation). While in Bayesian statistics, all quantities are probability distributions, so there is only the problem of inference.

TODO Look at this with more detail

- Deep Learning
 - Rich non-linear
 - Scalable learning
 - Easily composable
 - negative point
 - negative point
- Bayesian Reasoning
 - negative point: Rich non-linear
 - negative point: Scalable learning
 - negative point: Easily composable
 - positive point
 - positive point

Some additional examples of Bayesian reasoning combined with Deep learning are Density Estimation, Decision Making and Reinforcement Learning.

Deep and Hierarchical models can be decomposed into a sequence of conditional distributions in the form

$$p(z) = p(z_1|z_2)p(z_2|z_3) \dots p(z_{L-1}|z_L)p(z_L) \quad (2.5)$$

A list of different models

- Directed vs undirected
- observed vs unobserved variables
- parametric vs non-parametric

A list of learning principles (Statistical inference)

- Direct
 - Laplace approximation
 - Maximum likelihood
 - maximum a posteriori
 - variational inference
 - cavity methods
 - integr. nested laplace approximation
 - expectation maximisation
 - markov chain monte carlo
 - noise contrastive
 - sequential monte carlo

- Indirect (missing these ones)

Question about calibrated probabilities. In general in classification it is not but in healthcare is an example of application where this is important.

2.2 Inferential questions: Mon. 11:30–13:00

Evidence estimation

$$p(x) = \int p(x, z) dz \quad (2.6)$$

Moment computation

$$\mathbb{E}[f(x)|x] = \int f(z)p(z|x) dz \quad (2.7)$$

Identity trick to compute unobserved variables?

Integral problem

$$p(x) = \int p(x|z)p(z) dz \quad (2.8)$$

Probabilistic one

$$p(x) = \int p(x|z)p(z) \frac{q(z)}{q(z)} dz \quad (2.9)$$

Importance sampling: Monte carlo estimator

$$p(x) = \frac{1}{S} \sum_s w^{(s)} p(x|z^{(s)}) \quad (2.10)$$

$$w^{(s)} = \frac{p(z)}{q(z)} \dots \quad (2.11)$$

2.2.1 Hutchinson's Trick

Example with computing the trace of a matrix, KL between two gaussians, gradient of a log-determinant

The trace problem $Tr(A)$ with zero mean unit variance $\mathbb{E}[zz^T] = I$. Applying the identity trick we obtain

$$Tr(AI) = Tr(A \mathbb{E}[zz^T]) \quad (2.12)$$

With linear operations allow us to obtain

$$\mathbb{E}[Tr(Azz^T)] \quad (2.13)$$

And because of the trace property

$$\mathbb{E}[z^T Az]. \quad (2.14)$$

Then, with montecarlo methods the expected value is converted into a summation

2.2.2 Probability Flow Tricks

Given a distribution and sample

$$\hat{x} \sim p(x) \quad (2.15)$$

With a transformation

$$\hat{y} = g(\hat{x}; \theta) \quad (2.16)$$

Change of variables

$$p(y) = p(x) \left| \frac{dg}{dx} \right|^{-1} \quad (2.17)$$

Example: compute the

$$\log \det \left(\frac{\partial f(z)}{\partial z} \right) \quad (2.18)$$

This method is seen in the literature as *Normalising Flows* (see more in the blogpost)

2.2.3 Stochastic Optimization

Compute the gradient of the common problem

$$\nabla_{\phi} \mathbb{E}_{q_{\phi}(z)} [f_{\phi}(z)] = \nabla_{\phi} \int q_{\phi}(z) f_{\phi}(z) dz \quad (2.19)$$

With some reparametrisation tricks a distribution can be expressed as a transofmratio of other distributions

$$z \sim q_{\phi}(z) \quad (2.20)$$

Other names for these methods are samples, one-liners and change-of-variables

2.2.4 Pathwise Estimator

Also known as Unconscious statistician, stochastic backpropagation, perturbation analysis, reparametrisation trick, affine-independent inference.

When to use this trick

- Function f is differen
- Density q is known with a suitable transorm of a simpler base distribution: inverse CDF, location-scale transform, or other co-ordinate transform
- Easy to sample from base distribution

2.2.5 Log-derivative trick

Score function is the derivative of a log-likelihood function

$$\nabla_{\phi} \log q_{\phi}(z) = \quad (2.21)$$

2.2.6 Score-function estimator

Also known as Likelihood ratio method, reinforce and policy gradients, automated and black-box inference

We should use this method when

- Function is not differentiable, not analytical
- Distribution q is easy to sample from
- Density q is known and differentiable

2.2.7 Evidence Bounds

Integral problem

$$p(x) = \int p(x|z)p(z)dz \quad (2.22)$$

Proposal

$$p(x) = \int p(x|z)p(z) \frac{q(z)}{q(z)} dz \quad (2.23)$$

Importance weight

Jensen's inequality

Lower bound: evidence lower bound

$$\mathbb{E}_{q(z)} [\log p(x|z)] - KL[q(z)||p(z)] \quad (2.24)$$

2.2.8 Density Ratio Trick

The ratio of two densities can be computed using a classifier or using samples drawn from the two distributions. (TODO there is a type here, find out why)

$$\frac{p^*(x)}{q(x)} = \frac{p(y=1|x)}{p(y=-1|x)} \quad (2.25)$$

2.3 Unsupervised Learning Thu. 9:30–11:00

- Move beyond of associating inputs to outputs
- Generative models
- It allows to perform density estimation
 - Probabilistic models
 - High-dimensional data
 - Data distribution is targeted

Some examples of applications for generative models is the compression of images with high resolution GANs. Some artists used GANs to create pieces of video art. In an example the artist makes a video of a piece of cloth and gives the video as an input to a GAN that is trained with a particular set of images (eg. waves on the sea, or fire), then the GAN needs to generate new images that resemble the training data.

Types of generative models are:

- Fully-observed models
- Latent variable models (where there is a direction (from z to x ?)
- Undirected models: where there is no known direction from the observed variables X to the hidden variables Z
- <https://arxiv.org/abs/1202.3732>

Some points to consider when designing a generative model

- Data: binary, real-valued, nominal, strings, images
- Dependency: independent, sequential, temporal, spatial
- Representation: continuous or discrete
- Dimension: parametric or non-parametric
- Computational complexity
- Modelling capacity
- Bias, uncertainty, calibration
- Interpretability

2.3.1 Fully-observed models

These are models that operate on the observed data directly. It does not assume any hidden variables that may interact with our observed variables. As an example, Markov models assume a dependency with past events, depending on the number of dependencies with past events it has different orders of dependency.

$$x_1 \sim \text{Cat}(x_1|\pi) \quad (2.26)$$

$$x_2 \sim \text{Cat}(x_2|\pi(x_1)) \quad (2.27)$$

$$\dots \quad (2.28)$$

$$x_i \sim \text{Cat}(x_i|\pi(x_{<n})) \quad (2.29)$$

$$p(x) = \prod_i p(x_i|f(x_{<i}; \theta)) \quad (2.30)$$

Some of the properties of these models are

- Can directly encode how observed points are related
- Any type of data can be used
- ...

Directed and discrete: NADE, EoNADE... Directed and continuous: Normal means... Undirected and discrete... Undirected and continuous...

2.3.2 Latent variable models

These models introduce unobserved local random variables that represent a hidden cause. One of the most common assumptions is to assume some random hidden noise in the called Prescribed models. On the other hand, ...

An example of a prescribed model is a Deep Latent Gaussian Models in which all the hidden variables follow a Gaussian distribution that are connected one to each other and with dependencies to previous variables.

$$z_3 \sim \mathcal{N}(0, \mathbf{I}) \quad (2.31)$$

$$z_2 \sim \mathcal{N}(\text{some dependency on } z_3) \quad (2.32)$$

$$z_1 \sim \mathcal{N}(\text{some dependency on } z_3 \text{ and } z_2) \quad (2.33)$$

$$\dots \quad (2.34)$$

$$x_1 \sim \mathcal{N}(\text{some dependency on } z_i) \quad (2.35)$$

Some of the properties of Latent variable models are

- ...

Some dimensions in order to separate different latent models are, discrete vs continuous, deep vs direct/linear, and parametric vs non-parametric. Some examples are Buffet process, Sigmoid Belief Nets, Deep Gaussian processes, Hidden Markov Model, Sparse LVMS, Nonlinear factor Analysis.

2.3.3 Implicit Models

These are models that assume a random hidden variable that adds noise to the observed variables. One of the most common examples is to assume random Gaussian noise in a signal.

TODO the following equations need to be checked:

$$z \sim \mathcal{N}(\mu, \sigma^2) \quad (2.36)$$

$$x \sim f(z) \quad (2.37)$$

Some of the important properties are

- Easy to sample
- easy to compute expectations
- Can exploit on large-scale classifiers and Convolutional networks (I think this is the regularisation part?)

We can separate this models mostly on functions discrete time and diffusions in continuous time.

2.4 Model-Inference-Algorithm

We will understand VARIATIONAL Autoencoders and Generative ...

2.5 Variational Inference

The variational principle is a general family of methods for the approximation of a complicated density by a simpler class of densities. In most of the cases we can assess the similarity between our approximated distribution and the original one by using the Kullback–Leibler divergence.

IN variational inference there is always a variational bound that (Evidence Lower Bound, a.k.a. ELBO)

$$F(x, q) = \mathbb{E}_{q(z)} [\log p(x|z)] - KL[q(z)||p(z)] \quad (2.38)$$

2.6 Mean-Fields

These methods assume that the distribution is factorised (the hidden variables are independent). This means that we can compute their probability by multiplying every independent hidden variable probability.

The most expressive model with the true posterior would be $q^*(z|x) \propto p(x|z)p(z)$ (true posterior), in the least expressive side we have $q_{MF}(z|x) = \prod_z q(z_k)$ (fully-factorised model). There is a huge variety of models in between like Hidden Markov models, Autoregressive models, Gaussian processes?

2.7 Variational Optimisation

In the variational Expectation Maximisation consists on an alternating optimization for the variational parameters and the model parameters (VEM).

$$\phi \propto \nabla_{\phi} \text{ missing equation} \quad (2.39)$$

In the E-step instead of computing q with every sample of our dataset, we will simulate the answer with an Inference network $q(z/x)$ that will give us a sample $z \sim q(z/x)$. This may be an encoder or inversed. TODO need to see this previous paragraph with more detail.

2.8 Reinforcement learning as generative models

- An unknown likelihood
- not known analytical
- something more

Applying all the techniques seen in the three lectures it is possible to learn a policy learning that can be used in reinforcement learning.

Chapter 3

An Introduction to Gaussian Processes by Richard Turner

Example with a non-linear regression

Start with an example of a multivariate Gaussian

$$p(y|\Sigma) \propto \exp\left(\frac{-1}{2}y^T\Sigma^{-1}y\right) \quad (3.1)$$

If we condition the probabilities of one variable we obtain Gaussian shaped distributions as well

$$p(y_2|y_1, \Sigma) \propto \exp\left(\frac{-1}{2}y^T(y_2 - \mu_0)\Sigma^{-1}y(y_2 - \mu_0)\right) \quad (3.2)$$

TODO: revise previous equation

If we map any sample from the original distribution, it is possible to draw a line with the x-axis the different variables (variable index), and the y-axis the coordinates. If the variables are correlated, the lines should be quite horizontal.

In the examples it goes from 2 variables to 20, and show how the farther is the variable less correlated is to the first variable.

In the example, there is a covariance matrix with a Gaussian with fixed variance in the diagonal with the mean going from the first variable to the last one.

Then, if some of this variables are actual samples (fix values), then the sampling process should force the other variables to converge to these samples. (TODO: check if the following equation is an “l” or an “T”)

$$\Sigma(x_1, x_2) = \mathbf{K}(x_1, x_2) + l\sigma_y^2 \quad (3.3)$$

$$\mathbf{K}(x_1, x_2) = \sigma^2 e^{-\frac{1}{2l^2}(x_1 - x_2)^2} \quad (3.4)$$

I a non-parametric method (infinite parameters)

$$p(y|\theta) = \mathcal{N}(y; 0, \Sigma) \quad (3.5)$$

$$(3.6)$$

where σ^2 is the scale (vertical scale), and l is the horizontal-scale (this scales can be seen in the previous lineplot).

In a parametric model

$$y(x) = f(x; \theta) + \sigma_y \epsilon \quad (3.7)$$

$$\epsilon \sim \mathcal{N}(0, 1) \quad (3.8)$$

Definition of a Gaussian process: generalisation of a multivariate Gaussian distribution to infinitely many variables

A Gaussian distribution is fully specified by a mean vector, μ and a covariance matrix Σ .

$$f = (f_1, \dots, f_n) \sim \mathcal{N}(\mu, \Sigma), \text{ indices } i = 1, \dots, n \quad (3.9)$$

$$y(x) = f(x) + \epsilon\sigma_y \quad (3.10)$$

$$p(\epsilon) = \mathcal{N}(\epsilon; 0, 1) \quad (3.11)$$

place a GP prior over the non-linear function

$$p(f(x)|\theta) = GP(f(x); 0, K_\theta(x, x')) \quad (3.12)$$

$$\mathbf{K}(x, x') = \sigma^2 \exp\left(-\frac{1}{2l^2}(x - x')^2\right) \quad (3.13)$$

sum of a Gaussian variable into a GP is still a multivariate gaussian.

$$\text{missing equation} \quad (3.14)$$

The marginalisation property of Gaussian distributions is

$$p(y_1) = \int p(y_1, y_2) dy_2 \quad (3.15)$$

$$p(y_1, y_2) = \mathcal{N}(\text{[] missing}) \rightarrow p(y_1) = \mathcal{N}(y_1 : a, A) \quad (3.16)$$

How do we make a prediction

$$p(y_1|y_2) = \frac{p(y_1, y_2)}{p(y_2)} \quad (3.17)$$

$$\rightarrow p(y_1|y_2) = \mathcal{N}(y_1 : a + BC^{-1}(y_2 - b), A - BC^{-1}B^T) \quad (3.18)$$

Where y_1 are the predicted positions and y_2 are the samples

The predictive mean (first part of \mathcal{N} is linear in the data = Wy_2)

The predictive covariance (second part of the \mathcal{N}) can be interpreted as the predictive uncertainty = prior uncertainty A - the reduction in uncertainty $BC^{-1}B^T$.

What are the implications of the hyper-parameters?

- σ missing implications
- l missing implications

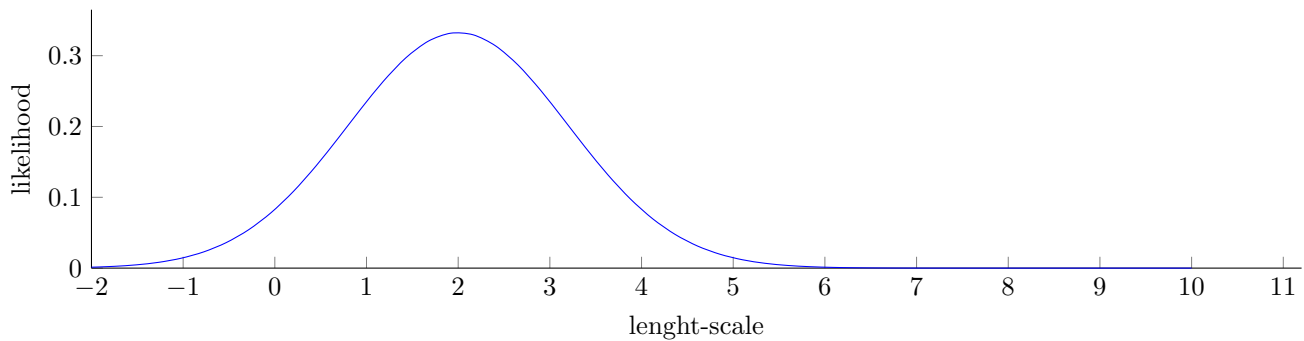
We can use the probability distributions to represent the plausibility of the hyper-parameters (uncertainty) given the data

$$\text{Bayes theorem with posterior probability of data given the parameters} \quad (3.19)$$

Shows an example of modifying the length-scale variable l and showing how the likelihood of the parameter value starts increasing and decreasing, showing a peak at lengthscale 2.

3.1 Covariance functions

Warning. Difficult to compare different models. If you try to compute the posterior probability of the different models given the data, it is hard to compute, it needs approximations and the results are really sensitive to the priors.

Figure 3.1: Likelihood of the data given the value of the length-scale l

3.1.1 References

- Gaussian Processes for Machine Learning, Rasmussen and Williams, 2006
- Gaussian Process Summer School, Neil Lawrence and colleagues
- Software
 - GPy: Gaussian processes in Python
 - GPflow: Gaussian Processes and tensorflow
 - GPML: Gaussian Processes in Matlab
 - GP Stan: Gaussian Processes in probabilistic programming

3.2 Using Gaussian Processes: Models, applications and connections 16:30–18:00

Question 1: Addition of two GPs

$$f(x) = f_1(x) + f_2(x) \quad (3.20)$$

$$f_1(x) \sim GP(0, \Sigma_1(x, x')) \quad (3.21)$$

$$f_2(x) \sim GP(0, \Sigma_2(x, x')) \quad (3.22)$$

The expected value is

$$m(x) = \mathbb{E}(f(x)) = \mathbb{E}(f_1(x) + f_2(x)) \quad (3.23)$$

$$= \mathbb{E}(f_1(x)) + \mathbb{E}(f_2(x)) \quad (3.24)$$

And the covariance is

$$\Sigma(x) = \mathbb{E}(f(x)f(x')) - \mathbb{E}(f(x))\mathbb{E}(f(x')) \quad (3.25)$$

$$= \mathbb{E}[(f_1(x) + f_2(x'))(f_1(x) + f_2(x')))] \quad (3.26)$$

$$= \mathbb{E}[f_1(x) + f_1(x')] + \mathbb{E}[(f_2(x) + f_2(x'))] + \mathbb{E}[f_1(x) + f_2(x')] + \mathbb{E}[(f_2(x) + f_1(x'))] \quad (3.27)$$

$$= \mathbb{E}[f_1(x) + f_1(x')] + \mathbb{E}[(f_2(x) + f_2(x'))] \quad (3.28)$$

More generally: GPs closed under linear transformation / combination:

- GP multiplied by a deterministic function = GP
- derivatives of GP = GP

- integral of a GP = GP
- convolution of a GP by a deterministic function = GP

Question 2: Random linear model

$$g(x) = mx + c \quad (3.29)$$

$$m \sim \mathcal{N}(0, \sigma_m^2) \quad (3.30)$$

$$c \sim \mathcal{N}(0, \sigma_c^2) \quad (3.31)$$

The expected value

$$m(x) = \mathbb{E}[g(x)] = \mathbb{E}(m)x + \mathbb{E}(c) = 0 + 0 \quad (3.32)$$

The covariance

$$\Sigma(x, x') = \sigma_m^2 xx' + \sigma_c^2 \quad (3.33)$$

$$= \mathbb{E}(g(x)g(x')) \quad (3.34)$$

$$= \mathbb{E}((mx + c)(mx' + c)) \quad (3.35)$$

$$= \mathbb{E}(m^2)xx' + \mathbb{E}(c^2) + \mathbb{E}(cm)x' + \mathbb{E}(mc)x \quad (3.36)$$

$$= \mathbb{E}(m^2)xx' + \mathbb{E}(c^2) + 0 + 0 \quad (3.37)$$

Question 3: random sinusoid model

$$h(x) = a \cos(wt) + b \sin(wt) \quad (3.38)$$

$$a \sim \text{Normal}(0, \sigma^2) \quad (3.39)$$

$$b \sim \text{Normal}(0, \sigma^2) \quad (3.40)$$

$$m(x) = 0 \quad (3.41)$$

$$\Sigma(x, x') = \sigma^2 \cos(w(x - x')) \quad (3.42)$$

Bochner's theorem: Any stationary covariance function can be written as

$$\Sigma(x - x') = \int \sigma^2(w) \cos(w(x - x')) dw \quad (3.43)$$

roughly, the function comprises and uncountably infinite sum of random sines and cosines

linear mappings $f(x) = Wx$	neural network	Gaussian Process mappings
linear regression	neural network regression	Gaussian Process regression
PCA or Factor analysis	variational auto-encoder (VAE)	Gaussian Process latent variable model
Gaussian auto-regressive (or Markov) model	neural auto-regressive density estimation (NADE)	Gaussian process auto-regressive model (GPAR)
linear Gaussian state space model (LGSSM)	recurrent neural latent variable model	Gaussian process state-space model (GP-SSM)

Strengths of Gaussian Processes

- Interpretable machine learning
- data-efficient machine learning
- decision making
- automated machine learning including probabilistic numerics

Weaknesses

- Large numbers of datapoints ($N \leq 10^5$)
- High-dimensional inputs spaces ($D \leq 10^2$)

In the speakers opinion, GPs are not good for large image recognition tasks, but for small tasks where it is crucial to get uncertainty estimates.

Example of an Interpretable auto-ML: the automatic statistician

Given some airline data it detects four underlying patterns, linearly increasing factor, a periodicity at every year, some increasing noise.

Shows an example with a video of an inverted pendulum and how in a small number of iterations (around 7?) a GP is learned.

3.2.1 Deep Gaussian Processes

Allowing non-parametric kernel spaces

$$y(x) = f(x) + \sigma_y \epsilon \quad (3.44)$$

$$\text{missing this part} \quad (3.45)$$

with a composition of GPs

$$y(x) = f(g(x)) + \sigma_y \epsilon \quad (3.46)$$

$$f(x) = GP(0, K_f(x, x')) \quad (3.47)$$

$$g(x) = GP(0, K_g(x, x')) \quad (3.48)$$

Deep GP may perform automatic kernel design

A Neural Network with one hidden layer and infinite number of units in the hidden layer is a Gaussian Process Nial 1996

A Neural Network with multiple hidden layers with infinite number of hidden units on each layer is also a GP (the variance on the weights need a specific variance σ^2/D Matthews et al. 2018. The specific variance is the reason why with finite number of units the regularisation needs to be readjusted and follows the values found by Matheews et al.

3.2.2 References

- Gaussian process latent variable models for visualisation of high dimensional data by Lawrence
- Local distance preservation in the GP-LVM through Back constraints
- the automatic statistician
- PILCO: a model-based and data-efficient approach to policy search

3.3 Large data and non-linear models Tue. 11:30–13:00

Shows a few examples of sounds generation using Gaussian Processes.

3.4 A brief history of gaussian process approximation

One of the first approaches was to modify the original samples in a way that the computational complexity of making exact inference was reduced from $O(n^3)$ to $O(n^2)$.

See following publications:

- Sparse Gaussian Processes using Pseudo-inputs
- Local and global sparse gaussian process approximations
- Sparse-posterior Gaussian Processes for general likelihoods
- Variational Learning of Inducing variables in sparse Gaussian Processes
- Fast Forward selection to speed up sparse Gaussian Process Regression

Some examples of Factor graphs

$$p(x_1, x_2, x_3) = g(x_1, x_2, x_3) \quad (3.49)$$

$$p(x_1, x_2, x_3) = g_1(x_1, x_2)g_2(x_2, x_3) \quad (3.50)$$

$$(3.51)$$

Where in the first one all the nodes are connected with one factor (square), while the second one x_1 is only connected to x_2 and this one to x_3 .

Shows an example of a multivariate gaussian where the covariance matrix has some numbers, and the inverse of the covariance matrix presents some zeros. By looking at the inverse of the covariance matrix we can see what nodes are independent given the rest (positions with the value zero).

$$\mathcal{N}(x, \mu, \Sigma) \propto \exp[-1/2(x - \mu)^T \Sigma^{-1}(x - \mu)] \quad (3.52)$$

$$= \exp[-1/2 \sum_{i,j} (x_i - \mu_i) \Sigma_{i,j}^{-1} (x_j - \mu_j)] \quad (3.53)$$

$$= \prod_{i,j} \exp[-1/2 (x_i - \mu_i) \Sigma_{i,j}^{-1} (x_j - \mu_j)] \quad (3.54)$$

$$= \prod_{i,j} g_{i,j}(x_i, x_j) \quad (3.55)$$

The previous equations show that when the inverse of the covariance is 0, the exponent value goes to one and the multiplicative factors are not affected.

The Kullback-Leibler divergence has the Gibb's inequality property, that means that $KL(p_1(z)||p_2(z)) \geq 0$ and has the equality at $p_1(z) = p_2(z)$. In order to prove the previous property apply the Jensen's inequality or differentiation. It is also Non-symmetric $KL(p_1(z)||p_2(z)) \neq KL(p_2(z)||p_1(z))$.

3.4.1 Fully independent training conditional (FITC) approximation

1. Augment model with $M < T$ pseudo data

$$p(f, u) = \mathcal{N}(f|u, \text{with mean } 0 \text{ and covariance } K) \quad (3.56)$$

2. Remove some of the dependencies (results in a simpler model)
3. Calibrate model (e.g. using KL divergence, many choices)

- This method introduces a parametric bottleneck into a non-parametric model
- Should we add extra pseudo-data when more data is available?

3.4.2 Variational free-energy method (VFE)

In this case we want to lower bound the likelihood

$$\mathcal{L}(\theta) = \log p(y|\theta) = \log \int df p(y, f|\theta) \quad (3.57)$$

$$= \log \int df p(y, f|\theta) \frac{q(f)}{q(f)} \geq \int df q(f) \log \frac{p(y, f|\theta)}{q(f)} = F(\theta) \quad (3.58)$$

$$= \int df q(f) \log \frac{p(f|y, \theta)p(y|\theta)}{q(f)} = \log p(y|\theta) - \mathbf{KL}(q(f)||p(f|y)) \quad (3.59)$$

$$F(\theta) = \int df q(f) \log \frac{p(y, f|\theta)}{\text{missing}} \quad (3.60)$$

At the end we get a lower bound of the likelihood.

$$F(\theta) = \int df q(f) \log \frac{p(y, f|\theta)}{p(f_{\neq u}|u)q(u)} \quad (3.61)$$

$$= \int df q(f) \log \frac{p(y|f, \theta)p(f_{\neq u}|u)p(u)}{p(f_{\neq u}|u)q(u)} \quad (3.62)$$

$$= \int df q(f) \log \frac{p(y|f, \theta)p(u)}{q(u)} \quad (3.63)$$

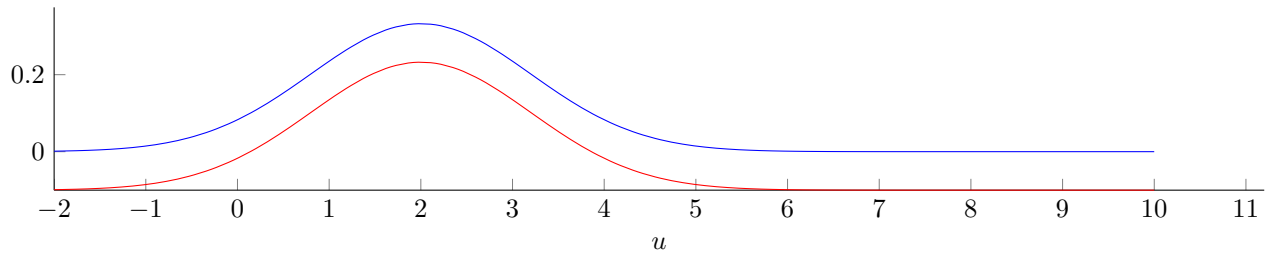


Figure 3.2: Where the blue line is the true likelihood and the red line is the approximated $q(u)$

They may reduce the time complexity from $O(M^3)$ to $O(NM^2)$, however N is proportional to M and in certain way is just a small improvement.

Some extra references about approximate inference in GPs:

- Sparse online gaussian processes
- A unifying view of sparse approximate gaussian process regression
- variational learning of inducing variables in sparse gaussian processes
- on sparse variational methods and the kullback-leibler divergence between stochastic processes
- A unifying framework for Gaussian Process Pseudo-Point approximation using Power Expectation propagation
- Streaming sparse gaussian process approximations
- Efficient deterministic approximate Bayesian Inference for Gaussian Process Models

and about Deep Gaussian processes:

- Deep gaussian processes for regression using approximate expectation propagation
- Doubly stochastic variational inference for deep gaussian processes

Chapter 4

Deep Learning by Rob Fergus

14:30–16:00

“Rob Fergus is an Associate Professor of Computer Science at the Courant Institute of Mathematical Sciences, New York University. He is also a Research Scientist at Facebook, working in their AI Research Group. He received a Masters in Electrical Engineering with Prof. Pietro Perona at Caltech, before completing a PhD with Prof. Andrew Zisserman at the University of Oxford in 2005. Before coming to NYU, he spent two years as a post-doc in the Computer Science and Artificial Intelligence Lab (CSAIL) at MIT, working with Prof. William Freeman. He has received several awards including a CVPR best paper prize, a Sloan Fellowship & NSF Career award and the IEEE Longuet-Higgins prize.” – Personal Page from New York University

4.1 Deep Supervised Learning

4.1.1 History of Neural Nets

The connections started around 40's with Hebb's work, McCulloch and Pitts and Rosenblatt 50's.

Second era started around 80's with Backpropagation to train multi-layered networks with the work of Rumelhart, Hinton and Williams. And the architecture of the Convolutional neural networks by Yann LeCun.

The era of Deep learning starts around 2011 with tasks focused on vision, speech understanding and natural language processing. The main ingredients were (1) supervised training of deep networks with (2) faster GPUs and (3) big labeled datasets.

4.1.2 Deep learning vs traditional approaches

The traditional approach consists on the hand-designed feature extraction that may be used by a simple classifier in order to predict the labels of a set of data. On the opposite hand, deep neural networks are focused on learning useful feature representations that improve the performance of the model on the dataset.

From 2010 to 2015 there is a clear decrease on the error rate in the ImageNet tasks by using more complex neural networks. The number of layers goes from 8 layers with the AlexNet in 2012 to 152 layers with the ResNet in 2016.

There are similar jumps in the performance of speech recognition systems. As an example, in 2015 Baidu's Deep Speech 2 system used 100 million parameters in 11 layers of a Recurrent Neural Network. It was trained in around 11.940 hours of English.

In the task of natural language processing and Machine Translation the work by Sutskever et al. and Cho et al. 2014. Tests of perplexity show a linear decrease from 2013?

One of the benefits of Deep learning is that although there is a huge computational cost during training, it is really lightweight during deployment. This means that small devices like smartphones are able to perform real-time predictions. One example is the detection of cars, road and pedestrians on self-driving cars.

Other interesting areas where deep learning is being applied are: astronomy, healthcare (e.g. skin cancer classification, or diabetic retinopathy).

4.1.3 Some issues with Deep Learning

(1) There is a missing theoretical understanding and performance guarantees. (2) Difficult to inspect the models. (3) Need lots of labeled data (4) We are still hand-designing the network architecture (instead of the features). There are some attempts to learn automatically the architecture (meta-learning) like Neural architecture search with reinforcement learning. (5) The most generic architectures (i.e. fully connected) does not seem to work, and it seems that their architecture is really dependent on the domain (e.g. for images exploit 2D)

4.1.4 Convolutional Neural Networks

First designed in the work of Yann LeCun (and previously Fukushima) by trying to mimic the visual system (of cats in Fukushima). It tries to exploit the spatial relation of pixels. In the work of Fukushima there was no back-propagation and the architecture could not be learned automatically. It was after back-propagation was applied into neural networks in 1986 that Yann LeCun could show a practical application by classifying handwritten digits from a dataset generated by some american post offices with Yann LeCun as an investigator (TODO rephrase previous sentence).

If the traditional activation functions were sigmoid shaped like the hyperbolic tangent or the logistic function, in more recent years activations with linear regions started to become more popular (e.g. ReLu).

Another useful trick is the Batch Normalisation, that has been empirically proved to improve the performance of CNNs. This method consists on normalizing the input features with 0 mean and 1 standard deviation on every mini-batch. There are two extra parameters (TODO missed these parameters).

4.1.5 Stochastic Gradient Descent

This is an iterative learning method that usually trains with mini-batches.

$$\text{missing} \quad (4.1)$$

Another training method is the AdaGrad

$$\theta_{t+1} = \theta_t - \text{missing} \quad (4.2)$$

RMSProp

$$\text{missing} \quad (4.3)$$

and ADAM

$$\text{missing} \quad (4.4)$$

4.1.6 Some practical debugging tips from M. Ranzato

- Train on small subset of data: the training error should go to zero
- Training diverges: learning rate may be too large (decrease learning rate), or the back-propagation is buggy because of numerical gradient issues.
- The parameters collapse, the loss is minimized but the training accuracy does not increase: check the loss function
- TODO some other tricks missing

4.1.7 Weights' initialization

With a simple example with a linear activation with 1 layer neural network we have

$$\text{Var}[y] = (n^{\text{in}} \text{Var}[w]) \text{Var}[x] \quad (4.5)$$

While in a multilayer network

$$\text{Var}[y] = \left(\prod_d n_d^{\text{in}} \text{Var}[w_d] \right) \text{Var}[x] \quad (4.6)$$

This value during the forward propagation will tend to explode with the depth d . And the backward propagation will make the gradients vanish.

$$\text{missing derivative of previous equation} \quad (4.7)$$

This makes the initialization really important for a good convergence of the weights into a good region within reasonable number of epochs (shows empirical comparison with good and bad initializations).

4.1.8 Deep Residual Learning

It allows some layers to bypass the next layer. This allows the network to skip layers if these are not necessary for the precision. Empirical experiments show a decrease from 7.4 error rate (with 34 layers) to 5.7 (with 152 layers).

Some additional information is described in “Tradeoffs of Large Scale Learning, Bottou & Bousquet, 2011”.

4.1.9 Weakly supervised pretraining

Some recent work on using weakly labeled images from the internet in order to pretrain deep neural networks.

Learning Visual Features from Large Weakly Supervised Data Joulin et al. (2015)

4.2 Unsupervised Learning

Learning without labels the intrinsic structure of the data. This is practically important as the real-world categories follow a Zipf’s law, meaning that lots of categories appear very few times.

Some of the benefits of unsupervised learning are that there is a vast amount of free data available, and it is potentially useful as a regularisation method.

The basic idea is to be able to build a model of $p(X)$ (density modeling) given just data $\{X\}$.

4.2.1 Self supervised learning

This method consists on using the input data as some part of the target data

$$y : X \rightarrow Y \quad (4.8)$$

$$x \rightarrow y(x) \quad (4.9)$$

and use the same techniques used in supervised learning to train a model

$$\arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n l(f_{\theta}(x_i), y(x_i)) \quad (4.10)$$

This approach usually requires some knowledge of the domain. Some examples are word2vec in which the word in the central position of a sentence is considered the target prediction.

4.2.2 Auto-Encoder

This method tries to find a hidden representation that keeps as much information of the input data as possible by reducing the reconstruction error. The network is divided between a decoder and an encoder part.

4.2.3 Variational Auto-Encoder

Similar to the simple autoencoder but makes

4.2.4 Generative Adversarial Networks

This consists in a decoder-only model but that uses an adversarial loss term using a discriminator D . Mini-max game between G and D.

4.2.5 Stacked Auto-Encoders

The Ladder Networks (Rasmus et al. 2015) adds a reconstruction constrain at every layer. In this way, there is a loss between every layer that can be minimized. This is like being able to generate hidden variable states apart from the input data.

4.2.6 Other approaches

Autoencoders, restricted / Deep Boltzmann Machines, ...

4.2.7 Application examples

Stochastic video generation (Denton and Fergus 2018) tries to predict the following frames from a video using encoders, decoders and Long Short Term Memories. The underlying method is a Variational Auto-Encoder.

An application example is shown with synthetically generated videos using MNIST dataset and two bouncing numbers in a closed square. The task of the model is to predict the position of the numbers in the future (pixel by pixel?). It is possible to sample from the learned distribution, possibly generating an un/certainty heat-map.

Training a generative model to add color to images from just their luminance levels.

Training a generative model to predict the word in the middle of a sentence.

Cut an image into different sections (3x3 squares?) and train a model to predict to which part of the image every pice belongs to. Unsupervised Learning of Visual Representations by Solving Jigsaw Puzzles Noroozi and Favaro (2016).

Automatic learning of images and audio from videos Owens et al. (2016). This method shows how it is possible to know the source of a sound from an image, some examples show people talking and rainfalls.

Map every image into a hypersphere Bojanowski and Joulin (2017), hopefully similar images will fall into nearby regions?

Unsupervised learning of visual representations using videos Wang and Gupta (2015)

4.3 Deep Learning Models Rob Fergus 9:30–11:00

Most of the deep learning models are exploiting some structural bias of the data. For example, images have a 2D correlation pixelwise, while speech recognition and natural language processing exploits a 1D time correlation.

4.3.1 Memory in Deep networks

Deep neural networks are not able to store memory to solve sequential problems. Some networks are able to store implicit internal memory on their connections like Recurrent Neural Networks and Long Short Term Memory networks.

Implicit internal memory

Recurrent neural networks have the computational and the memory integrated in their weights. However, one of the main problems of RNNs is how to prevent the network from forgetting during the training. In order to solve that problem Mikolov et al. 2014 separated the state into a fast and slow changing part. Other approaches are by using gated units for the internal state like Long-short term memory (LSTM) (see Hochreiter & Schmidhuber 1997, Graves, 2013), or Simplified minimal gated unit variations for RNNs. Also GRU light-weight version from Cho et al 2014.

RNN search: attention in machine translation Bahdanau et al. (2014) investigates an encoder and decoder model in a RNN, similar method was used later for image caption generation with attention Xu et al. (2015), this last method allows the network to be queried and create heatmaps on the original images to correlate the generated words with their corresponding patch.

External Global Memory

Some work tries to split the computational memory from the problem solving memory aspect. These two parts are called the *memory module* and the *controller module*. The controller is able to write and read to the memory

module. The memory needs some addressing mechanism, and can be soft or hard. The benefit of soft addressing is that it can be trained by backpropagation, while the hard addressing is not differentiable.

The end-to-end memory network (MemN2N) Sukhbaatar et al. (2015) is an architecture. The memory module is able to store some input vectors and hidden states, the controller module is able to search in the memory for similar patterns and obtain the corresponding associated output. This can be a new hidden state that will be combined with the current hidden state of the controller.

1. Embed each word (word to vector?)
2. Sum embedding vectors ($v_{Sam} + v_{drops} + v_{apple} = m_i$)
3. It generates one memory per sentence
4. The controller takes the vector encoding of the question and makes a dots product with the memory sentences and applies a softmax.
5. This gives weights to each sentence it computes a weighted sum of the sentences' representations.
6. The representation is given to the controller and generate an answer from it.

Another example is the Neural Turing Machine Graves et al. (2014) in which the authors design an end-to-end differentiable architecture that may be trained by backpropagation to solve tasks involving memory. Some example applications are coping, sorting and associative recall of inputs with outputs.

4.4 Deep Nets for sets

There are problems where there is permutation invariance, dynamic sizing, single output, or output for each element.

In the Communication Neural Network (CommNet) Sukhbaatar et al. (2016) the inputs and outputs are sets, another example is the DeepSets Zaheer et al. (2017). The CommNet is a special case of a Graph NN in which a set can be represented as a complete graph.

Bibliography

- Bahdanau, D., Cho, K. and Bengio, Y. (2014). Neural machine translation by jointly learning to align and translate, *arXiv preprint arXiv:1409.0473*.
- Bojanowski, P. and Joulin, A. (2017). Unsupervised learning by predicting noise, *arXiv preprint arXiv:1704.05310*.
- Graves, A., Wayne, G. and Danihelka, I. (2014). Neural Turing machines, *arXiv preprint arXiv:1410.5401*.
- Joulin, A., van der Maaten, L., Jabri, A. and Vasilache, N. (2015). Learning visual features from large weakly supervised data, *CoRR* **abs/1511.02251**.
URL: <http://arxiv.org/abs/1511.02251>
- Noroozi, M. and Favaro, P. (2016). Unsupervised learning of visual representations by solving jigsaw puzzles, *CoRR* **abs/1603.09246**.
URL: <http://arxiv.org/abs/1603.09246>
- Owens, A., Wu, J., McDermott, J. H., Freeman, W. T. and Torralba, A. (2016). Ambient sound provides supervision for visual learning, *CoRR* **abs/1608.07017**.
URL: <http://arxiv.org/abs/1608.07017>
- Sukhbaatar, S., Fergus, R. et al. (2016). Learning multiagent communication with backpropagation, *Advances in Neural Information Processing Systems*, pp. 2244–2252.
- Sukhbaatar, S., Szlam, A., Weston, J. and Fergus, R. (2015). Weakly supervised memory networks, *CoRR* **abs/1503.08895**.
URL: <http://arxiv.org/abs/1503.08895>

- Wang, X. and Gupta, A. (2015). Unsupervised learning of visual representations using videos, *Proceedings of the IEEE International Conference on Computer Vision*, pp. 2794–2802.
- Xu, K., Ba, J., Kiros, R., Cho, K., Courville, A., Salakhudinov, R., Zemel, R. and Bengio, Y. (2015). Show, attend and tell: Neural image caption generation with visual attention, *International conference on machine learning*, pp. 2048–2057.
- Zaheer, M., Kottur, S., Ravanbakhsh, S., Poczos, B., Salakhutdinov, R. R. and Smola, A. J. (2017). Deep sets, *Advances in Neural Information Processing Systems*, pp. 3391–3401.

Chapter 5

Statistical machine learning and convex optimization by Francis Bach

5.1 Introduction 14:30–18:00

We will use n for the number of observations and d for the dimension size. When working with big data we will seek ideally for a running-time complexity of $O(dn)$.

5.2 Classical methods for convex optimization

In supervised learning we have a set of observations $(x_i, y_i) \in X, Y, i = 1, \dots, n$, i.i.d. (this assumption is almost never true).

For regression $y \in \mathbb{R}$ and prediction $\hat{y} = \theta^T \Phi(x)$, for classification $y \in \{-1, 1\}$ and the prediction is $\hat{y} = \text{sign}(\theta^T \Phi(x))$.

Empirical risk minimization to find a $\hat{\theta}$ solution to

$$\arg \min_{\theta \in \mathbb{R}^d} 1/n \sum_{i=1}^n l(y_i, \theta^T \Phi(x_i)) + \mu \Omega(\theta) \quad (5.1)$$

the training cost (or empirical risk) is

$$\hat{f}(\theta) = 1/n \sum_{i=1}^n l(y_i, \theta^T \Phi(x_i)) \quad (5.2)$$

and the testing cost or expected risk is

$$\hat{f}(\theta) = \mathbb{E}_{(x,y)} l(y_i, \theta^T \Phi(x_i)) \quad (5.3)$$

This imposes two fundamental questions: (1) how are we finding the optimal set of parameters θ and (2) how to analyse $\hat{\theta}$.

The loss for a single observation is $f_i(\theta) = l(y_i, \dots)$ TODO missing

Jensen inequality says that $g(\mathbb{E}(\theta)) \leq \mathbb{E} g(\theta)$

The global definition of convexity if we assume differentiability is

$$\forall \theta_1, \theta_2, g(\theta_1) \leq g(\theta_2) + g'(\theta_2)^T (\theta_1 - \theta_2) \quad (5.4)$$

With twice differentiable functions $\forall \theta, g''(\theta) \succeq 0$ (positive semi-definite Hessians).

We ideally want a convex problem because the local minimum is the same as the global minimum, with an optimal condition in $g'(\theta) = 0$. See also the convex duality and recognizing convex problems in (Boyd and Vandenberghe; 2004).

5.2.1 Lipschitz continuity

Bounded gradients of g (lipschitz-continuity) if the function g is convex, differentiable and has a

5.2.2 Smoothness and strong convexity

Theorem 5.1. *A function $g : \mathbb{R}^d \rightarrow \mathbb{R}$ is L -smooth if and only if it is differentiable and its gradient is L -Lipschitz-continuous*

$$\forall \theta_1, \theta_2 \in \mathbb{R}^d, \|g'(\theta_1) - g'(\theta_2)\|_2 \leq L\|\theta_1 - \theta_2\|_2 \quad (5.5)$$

Adding a regularization $\mu/2\|\theta\|^2$ with a value of μ on the order of $1/n$.

5.2.3 Analysis of empirical risk minimization

most important slide from the first part is summarized in slide 47

5.2.4 Accelerated gradient methods (Nesterov, 1983)

Algorithm

$$\theta_t = \eta_{t-1} - \frac{1}{L}g'(\eta_{t-1}) \quad (5.6)$$

$$\eta_t = \theta_t + \frac{t-1}{t+2}(\theta_t - \theta_{t-1}) \quad (5.7)$$

with bound

$$g(\theta_t) - g(\theta_*) \leq \frac{2L\|\theta_0 - \theta_*\|^2}{(t+1)^2} \quad (5.8)$$

5.2.5 Optimization from sparsity-inducing norms

See Bach, Jenatton, Mairal, and Obozinski, 2012b (Bach et al.; 2012).

Newton method

Minimize the second-order Taylor expansion

5.2.6 Summary about minimization of convex functions

- Gradient descent: $\theta_t = \theta_{t-1} - \gamma_t g'(\theta_{t-1})$
 - $O(1/\sqrt{t})$ convergence rate for non-smooth convex functions
 - $O(1/t)$ convergence rate for smooth convex functions
 - $O(e^{-pt})$ convergence rate for strongly smooth convex functions
- Newton method: $\theta_t = \theta_{t-1} - g''(\theta_{t-1})^{-1}g'(\theta_{t-1})$
 - $O(e^{-p2^t})$ convergence rate
- Key insights from Bottou and Bousquet (2008)
 - In machine learning, it is not necessary to optimize below statistical error
 - In machine learning the cost functions are averages
 - Testing errors are more important than training errors

5.3 Convex stochastic approximation

There are known global minimax? rates of convergence for non-smooth problems (Nemirovsky and Yudin, 1983; Agarwal et al., 2012).

- Least-squares regression is easy to analyze, and has an explicit relationship to bias/variance trade-offs (see Défossez and Bach (2015); Dieuleveut et al. (2016)).
- Many important loss functions are not quadratic (see Bach and Moulines (2013)).

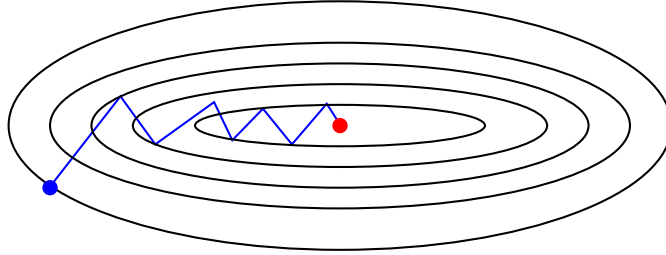


Figure 5.1: Batch Gradient Descent: it needs to decrease the error on every step, if not, then the parameters are not right

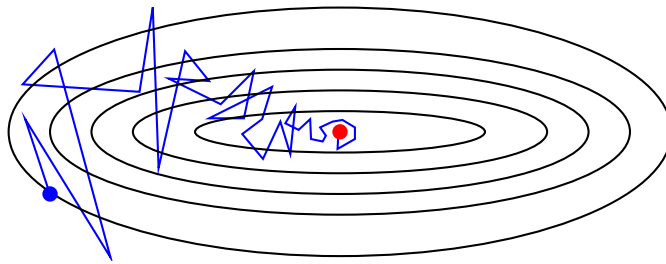


Figure 5.2: Stochastic Gradient Descent: the error may increase occasionally

5.4 Summary of rates of convergence

Given the problem parameters

- D : diameter of the domain
- B Lipschitz-constant
- L smoothness constant
- μ strong convexity constant

5.5 Conclusions

- Statistics with or without optimization
 - Significance of mixing algorithms with analysis
 - Benefits of mixing algorithms with analysis
- Open problems
 - Non-parametric stochastic approximation (Dieuleveut and Bach, 2014)

	convex	strongly convex
nonsmooth	deterministic: BD/\sqrt{t} stochastic: BD/\sqrt{n}	deterministic: $B^2/(t\mu)$ stochastic: $B^2/(n\mu)$
smooth	deterministic: LD^2/t^2 stochastic: LD^2/\sqrt{n} finite sum: n/t	deterministic: $\exp(-t\sqrt{\mu/L})$ stochastic: $L/(n\mu)$ finite sum: $\exp(-t/(n + L/\mu))$
quadratic	deterministic: LD^2/t^2 stochastic: $d/n + LD^2/n$	deterministic: $\exp(-t\sqrt{\mu/L})$ stochastic: $d/n + LD^2/n$

Table 5.1: Summary of rates of convergence

- Characterization of implicit regularization of online methods item Structured prediction
- Going beyond a single pass over the data (/testint performance)
- Parallel and distributed optimization
- Non-convex optimization (Reddi et al., 2016)

Bibliography

- Bach, F., Jenatton, R., Mairal, J., Obozinski, G., et al. (2012). Optimization with sparsity-inducing penalties. *Foundations and Trends® in Machine Learning*, 4(1):1–106.
- Boyd, S. and Vandenberghe, L. (2004). *Convex optimization*. Cambridge university press.

Chapter 6

Supervised Learning and Text Classification by Kyunghyun Cho

6.1 Introduction to supervised learning with ANNs Wed. 11:30–13:00

An overview on supervised learning, in which the input is described as a validation set D_{val} and a test set D_{test}

At the end we need to choose the hypothesis that best adjusts to our problem between the availables. Examples of hyperparameters sets are in SVMs the regularization parameters and the kernel, similarly with Gaussian Processes with the kernel and the parameters σ^2 and *length-scale*.

In Artificial Neural Networks the hypothesis set is the set of architectures, and hyperparameters. The architecture is commonly an directed acyclic graph (DAG) with parameters, inputs, outputs and compute nodes (functions that are often differentiable).

An example of architecture is a logistic regression where

$$p_{\theta}(y = 1|x) = \sigma(w^T x + b) = \frac{1}{1 + \exp(-w^T x - b)} \quad (6.1)$$

or a 3rd-order polynomial function.

The inference is done by forward propagation of the input to the output through all the hidden layers.

Supervised learning tries to find a function $f_{\theta}(x)$, while in the case of Neural Networks we can usually interpret it as computing the conditional probabilities given the input $p(y = ' | x)$. This is achieved by computing any arbitrary output values from a network, then exponentiating every individual prediction and dividing them by their sum (soft-max function).

The objective is that the training data is maximally maximized, by ensuring that the maximum individual probabilities is maximized at the same time.

$$\arg \max_{\theta} \log p_{\theta}(D) = \arg \max_{\theta} \sum_{n=1}^N \log p_{\theta}(y_n | x_n) \quad (6.2)$$

We can also use the log-likelihood

$$\text{Missing equation} \quad (6.3)$$

6.1.1 Loss minimization

In order to minimize the loss it is necessary to use optimization techniques, in the case of Artificial Neural Networks (ANNs) is by using backpropagation to compute the partial derivatives of the parameters with respect to the input using the chain rule of derivatives

$$\frac{\partial(f \circ g)}{\partial x} = \frac{\partial f}{\partial g} \frac{\partial g}{\partial x} \quad (6.4)$$

This differentiation is done automatically by autograd which will implement the Jacobian-vector product of each P node.

By doing backpropagation we obtain the gradient of the loss with respect to every parameter θ . Instead of computing the gradient for the full dataset, it is common to use mini-batches with a method called Stochastic Gradient Descent.

1. Random subset of M training examples
2. Compute the minibatch gradient

$$\nabla L \sim 1/N' \sum_{n=1}^{N'} \text{missing part} \quad (6.5)$$

3. Update the parameters missing equation
4. repeat until the validation loss stops improving

However, this method will find a local minima in the training set, in order not to overfit to the training data, a validation set is used to do an early stop. This is one of the most important parts of the training as we want the minima to be as near as possible to the dataset near.

6.2 Text classification

In text classification the input are sentences or paragraphs and the output is a category to which the input belongs to (commonly a fixed number of C categories).

Some of the particularities of using sentences as inputs is that they are of variable size $X = (x_1, x_2, \dots, x_T)$ where every x_t is a token from a vocabulary V . This is done by automatically (or manually) splitting all the individual words and creating an index of words that will form our vocabulary V . Then every sentence is encoded by a sequence of integers. In this case we are not tokenizing the words first by extracting any meaning of the word, in this sense the words “cat” and “cats” have completely independent indices.

At the end the words are encoded as one-hot-encoding with a bit to 1 for the corresponding index. This will be given to the ANN and will be forward propagated to obtain a hidden representation e_i . This can be interpreted as a table lookup from a word to a hidden embedding (eg. a fixed matrix W of dimension $\#tokens \times \#arbitrary\ dimension$ that is multiplied by the binary token).

The representation of a sentence is a continuous bag-of-words, this means that we ignore the order of the words in the sentence, and average the corresponding token vectors. The method has been proven to be very useful in the works Iyyer2016, cho2017, and in FastText by Bojanowski2017.

6.2.1 How to represent a sentence

When the order of the words is important, it is possible to encode every set of words. For example, Relation Network: skip bigrams consider all the possible pairs of tokens and averages all the relationship vectors. The relations between words can be encoded depending on the problem, we could consider that the order of the bigrams is not important, only interested on bigrams of contiguous words, all the possible pairs, or skip some words.

With the previous idea is possible to use Convolutional Neural Networks in order to convolve the “look up table” through the input Kim (2014).

The CNN can also represent the bigram representations, and it is possible to learn a weight vector that will choose how many words are important. Look at the relation between Relational Networks and Convolutional Neural Networks.

Self-attention 11:30–13:00

Self-attention is a generalization of a CNN and RN that is able to capture long-range dependencies with a single layer. It can also ignore irrelevant long-term dependencies. Also mention generalization with multi-head and multi-hop attention.

Using a Recurrent Neural Network we can create an online representation of a sentence by reading every word and storing their representation in the recursive hidden representation. This allows a cost of $O(T)$ instead of the $O(T^2)$ necessary to do all the word pairs.

The representation that is generated by the RNN can encode the representation of a region of the text given its context. This can then be feeded to the previously seen attention model that can learn the weighted sum of the context.

In “Fully character-level neural machine translation without explicit segmentation” Lee et al. (2016) the authors stack a RNN on top of CNNs.

6.3 Natural Language Models

6.3.1 Autoregressive language modelling

The autoregressive sequence modelling assumes that the past tokens influence the current token as

$$p(X) = p(x_1)p(x_2|x_1) \cdots p(x_T|x_1, \dots, x_{T-1}) \quad (6.6)$$

this holds true given the conditional distribution assumption.

With this method, an unsupervised method is transformed in T smaller supervised problems.

One think to keep in mind from a question that was asked is that although the marginalisation of $p(X)$ should sum to one, in real cases with RNNs this is not always true. Possibly because of the parametrization.

The autoregressive sequence modelling can be represented as

$$p(X) = \prod_{t=1}^T p(x_t|x_{<t}) \quad (6.7)$$

In order to score a sentence we can compare the output of the Autoregressive model for several sentences and use softmax to obtain “probabilities” (values between 0 and 1 that sum to one) for each sentence.

6.3.2 N-Gram language models

Before ANNs were used, the idea about ussing the conditional probabilities was already used on a smaller scale with N-gram models. In this case the N needs to be decided in advance.

$$p(x|x_{-N}, x_{-N+1}, \dots, x_{-1}) = \frac{p(x_{-N}, x_{-N+1}, \dots, x_{-1}, x)}{\sum_{x \in V} p(x_{-N}, x_{-N+1}, \dots, x_{-1}, x)} \quad (6.8)$$

The process then is to get the dataset and count the frequencies of every occurrence of the tokens \dots, x_{N-1} followed by all the possible tokens x .

There are two main issues that arise by ussing frequentist N-Grams

1. Data sparsity: lack of generalization
 - If using a 3-gram model you find 3 words that never happened again, the product of probabilities will be zero independently on the rest.
 - One possible solution is to use smoothing by adding a small constant
 - Another solution is to try with all $n \in \{N, \dots, 1\}$
2. Inability to capture long-term dependencies
 - By choosing a fixed N we may lose long term dependencies.

6.3.3 Neural N-Gram Language Model

Bengio et al. (2003) solved some of the previous issues by using the hidden representation of a Neural Network instead of the tokens. In the continuous vector space the similar tokens or phrases are nearby.

Some other work on the same direction is Mikolov et al. (2013), , Pennington et al. (2014), Le and Mikolov (2014).

An example of this application is the generalization of sentences that never happened by realizing that some words share some similarities (eg. numbers). An example was shown in which sentences with the 2-grams “three teams”, “four teams” and “four groups” are able to generalise to the bigram “three groups” by realizing that three and four share a similar continuous space, and that before groups there could be a number.

In practice

1. Collect TODO missing steps

6.3.4 Convolutional Language Models

Convolutional Neural Networks allow to extend the context size by applying the convolution through larger parts of the text (see kalchbrenner et al, 2015 and Dauphin et al 2016, ByteNet by Kalchbrenner et al. (2016), PixelCNN, WaveNet, ...)

Gated convolutional language model by Dauphin 2016 Dauphin et al. (2016)

6.3.5 CBoW Language Models (infinite context)

The idea is similar to the LM of using averages instead of concatenation.

6.3.6 Recurrent Language Models

An RNN can summarize all the tokens seen until x into a continuous vector representation Mikolov et al, 2010 Mikolov et al. (2010).

6.3.7 Recurrent Memory Networks

The work of Tran et al., 2016 Tran et al. (2016) combines RNNs to compress the context into a continuous vector representation with the attention model that learns the weighting of the context.

6.4 Recurrent Networks and Backpropagation

Consider the full path from a parameter θ to the loss l . The backpropagation consists on computing the gradient of the l with respect to the previous node and multiply by the Jacobian matrix of every step back to the parameter θ .

$$\mathbf{Jac}_{h^{t+1}}^{h^t} = W^T \text{diag}(\tanh'(a^t)) \quad (6.9)$$

Because the Jaciobians are multiplied by every backpropagation step, this means that

- If $W > 1$ (the upperbound of) the norm blows up, exploding gradient
- If $W < 1$ (the upperbound of) the norm shrinks to zero, vanishing gradient

$$\| \prod_{t'=t} \| \leq \| \mathbf{Jac} \| \text{TODO missing part} \quad (6.10)$$

6.4.1 Gated recurrent units

These type of networks allow to skip some of the paths in order to avoid the exploding or vanishing gradient.

- Adaptive shortcut:
- Candidate update + pruning
- Update gate: $u_t = \sigma(W_u h_{t-1} + U \dots$
- reset gate $r_t = \sigma(W_r h_{t-1} + U_r x + b_r)$

6.4.2 Lessons from GRU/LSTM

- Credit assignment over a long path of computation is difficult
- Adaptive shorcut or skip-connection helps avoid credit dilution
- Gates are an effective way to focus credit assignment

6.5 Neural Machine Translation 14:30–16:00

6.5.1 History of machine translation

The original idea was to get a text from one language, (1) perform a morphological analysis, (2) a syntactic analysis, (3) semantic analysis, (4) a semantic composition and obtain an interlingua text that can be transformed back to any other language Borr, Hovy and Levin 2006 **TODO missing citation**.

Allen 1987 *iecc icnn*, a brief resurrection of Neural Networks in 1997 by Castano and Casacuberta 1997 Castano et al. (1997), then in 2006 Schwenk 2006 Schwenk et al. (2006) as a filter source to SMT to ANN to target sentence, then Devlin 2014 from source to SMT + ANN to target sentence, then source to ANN to target sentence.

6.5.2 Encoding: Token representation

First it is necessary to build a source and target vocabulary of unique tokens (for each language). Then transform the text into the set of tokens. Then encode the token sentences into sentence representation vectors, being careful not to compress the sentences into small vectors that may lose useful information.

6.5.3 Decoding: conditional language modeling

Using autoregressive networks we are interested in predicting the posterior probability

$$p(Y|X) = \prod_{t=1}^T p(y_t|y_{t-1}, X) \quad (6.11)$$

Look at the RNN Neural Machine Translation by Bahdanau et al., 2015 Bahdanau et al. (2014). The model uses the target sentence (and what has been translated until the current moment) in order to generate the following prediction.

1. Encode: read the entire source sentence to know what to translate
2. Attention:
3. Decode:
4. Repeat 2 to 3 until the end-of-sentence (token) is achieved

This method achieved performance as good as the current state-of-the-art alternative at the moment phrase-based machine translation (PBMT).

At translation time every predicted word of the sentence had an associated vector of weights that indicted the source words involved, although the model was trained only with pairs of text without any extra supervision. We should consider that every token of the source sentence is at the same time associated with a context in its language (see Jean et al. (2015)).

6.5.4 In practice

Available frameworks

- Nematus Sennrich et al. (2017)
- OpenNMT-py Klein et al. (2017)
- FairSeq Gehring et al. (2017)
- Sockeye Hieber et al. (2017)

6.6 Current and ongoing projects

Multilingual translation, real-time translation, and character level translation.

6.6.1 Multilingual translation

A common approach has been to use a pivot language to translate languages without lots of examples. For example, translate Korean to Japanese and then to English as the corpus between them is larger than the Korean to English.

With this idea, there has been some work to generate a pivot common language that is automatically learned in an ANN. For example, an encoder/decoder approach in Firat et al. (2016a,b) creates one encoder per source and decoder per target, and a model is learned for every pair. Later Johnson et al, 2016, Ha et al, 2016, Lee et al 2017, and Gu et al, 2018 work on an Universal ...

A current limitation is these methods drastically depend on the different amount of available data in each language. The models start ignoring the less frequent language. However, if given the same proportions of samples for every language they may overfit to the less frequent one because of the multiple epochs run on one language compared to the other.

Some work to solve this problem is being done with Meta-Learning MAML in Finn et al. (2017). The difference with multitask learning is that ...

Another idea is to create the unsupervised lookup tables to convert word tokens into continuous vectors using big corpus of different languages. In this manner, words with similar meanings will fall into similar regions (see Artetxe et al. (2017)). In this case, the overfitting is less probable to happen as the lookup table is only learning a mapping of tokens to a continuous space, but not the translation. Then, the translation model is trained on top of it.

6.6.2 Real-Time Translation (learning to decode)

learning to translate in real-time with neural machine translation Gu et al. (2016)

In order to generate the best translation it is possible to generate several and then select the most probable one. However, occasionally it may happen that after several words are translated the topic changes and the best translation may be a different one. In order to solve this the Beam Search keeps track of parallel best translations and is able to switch to another one at any point if necessary.

Exploiting the hidden activation

In a Deep Neural Network huge information in the hidden layers is usually discarded in order to predict a class; obtaining at the end a binary prediction. However, the hidden information has rich information about the original input that can be exploited.

By using the hidden information the performance of several methods was increased Gu et al. (2016).

Bibliography

- Artetxe, M., Labaka, G., Agirre, E., and Cho, K. (2017). Unsupervised neural machine translation. *arXiv preprint arXiv:1710.11041*.
- Bahdanau, D., Cho, K., and Bengio, Y. (2014). Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*.
- Bengio, Y., Ducharme, R., Vincent, P., and Jauvin, C. (2003). A neural probabilistic language model. *Journal of machine learning research*, 3(Feb):1137–1155.
- Castano, M. A., Casacuberta, F., and Vidal, E. (1997). Machine translation using neural networks and finite-state models. *Theoretical and Methodological Issues in Machine Translation (TMI)*, pages 160–167.
- Dauphin, Y. N., Fan, A., Auli, M., and Grangier, D. (2016). Language modeling with gated convolutional networks. *arXiv preprint arXiv:1612.08083*.
- Finn, C., Abbeel, P., and Levine, S. (2017). Model-agnostic meta-learning for fast adaptation of deep networks. *arXiv preprint arXiv:1703.03400*.
- Firat, O., Cho, K., and Bengio, Y. (2016a). Multi-way, multilingual neural machine translation with a shared attention mechanism. *arXiv preprint arXiv:1601.01073*.

- Firat, O., Sankaran, B., Al-Onaizan, Y., Vural, F. T. Y., and Cho, K. (2016b). Zero-resource translation with multi-lingual neural machine translation. *arXiv preprint arXiv:1606.04164*.
- Gehring, J., Auli, M., Grangier, D., Yarats, D., and Dauphin, Y. N. (2017). Convolutional Sequence to Sequence Learning. *ArXiv e-prints*.
- Gu, J., Neubig, G., Cho, K., and Li, V. O. (2016). Learning to translate in real-time with neural machine translation. *arXiv preprint arXiv:1610.00388*.
- Hieber, F., Domhan, T., Denkowski, M., Vilar, D., Sokolov, A., Clifton, A., and Post, M. (2017). Sockeye: A Toolkit for Neural Machine Translation. *arXiv preprint arXiv:1712.05690*.
- Jean, S., Firat, O., Cho, K., Memisevic, R., and Bengio, Y. (2015). Montreal neural machine translation systems for wmt’15. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 134–140.
- Kalchbrenner, N., Espeholt, L., Simonyan, K., Oord, A. v. d., Graves, A., and Kavukcuoglu, K. (2016). Neural machine translation in linear time. *arXiv preprint arXiv:1610.10099*.
- Kim, Y. (2014). Convolutional neural networks for sentence classification. *arXiv preprint arXiv:1408.5882*.
- Klein, G., Kim, Y., Deng, Y., Senellart, J., and Rush, A. M. (2017). OpenNMT: Open-source toolkit for neural machine translation. In *Proc. ACL*.
- Le, Q. and Mikolov, T. (2014). Distributed representations of sentences and documents. In *International Conference on Machine Learning*, pages 1188–1196.
- Lee, J., Cho, K., and Hofmann, T. (2016). Fully character-level neural machine translation without explicit segmentation. *arXiv preprint arXiv:1610.03017*.
- Mikolov, T., Karafiát, M., Burget, L., Černocký, J., and Khudanpur, S. (2010). Recurrent neural network based language model. In *Eleventh Annual Conference of the International Speech Communication Association*.
- Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S., and Dean, J. (2013). Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pages 3111–3119.
- Pennington, J., Socher, R., and Manning, C. (2014). Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543.
- Schwenk, H., Dchelotte, D., and Gauvain, J.-L. (2006). Continuous space language models for statistical machine translation. In *Proceedings of the COLING/ACL on Main conference poster sessions*, pages 723–730. Association for Computational Linguistics.
- Sennrich, R., Firat, O., Cho, K., Birch, A., Haddow, B., Hitschler, J., Junczys-Dowmunt, M., Läubli, S., Barone, A. V. M., Mokry, J., et al. (2017). Nematus: a toolkit for neural machine translation. *arXiv preprint arXiv:1703.04357*.
- Tran, K., Bisazza, A., and Monz, C. (2016). Recurrent memory networks for language modeling. *arXiv preprint arXiv:1601.01272*.

Chapter 7

Causality by Joris Mooij

There are many questions in science that are casual (eg. climatology, healthcare, ...).

7.1 Introduction

Causation is not correlation (gives example with chocolate consumption being correlated with number of nobel prizes in different countries, while there is no actual correlation between both).

In order to represent causal relations we can use *causal graphs* (directed graphs) in wich nodes are variables X_n from (a vocabulary?) V . While directed edges indicate that the first variable causes directly another variable respect to V . An example of a cyclic graph where every adjacent variable is directly connected to its neighbours is the encoding of standing domino pieces.

It is possible to modify the *causal graph* with an *intervention*, in the example the domino piece X_2 is glued to the floor, removing the direct connections from other pieces to X_2 , but possibly keeping a connection from X_2 to the adjacent (in case an external user can still force X_2 to fall).

A *perfect* (“surgical”, “atomic”) *intervention* on a set of variables $X \subseteq V$, is an external enforced change of the system (eg. the previous example).

A *confounder* is a latent common cause:

H is a confounder of X and Y if:

1. H causes X directly w.r.t. $\{X, Y, H\}$
2. H causes Y directly w.r.t. $\{X, Y, H\}$

We will denote latent confounders by *bidirected edges* in a causal graph.

7.2 Defining causality in terms of probabilities

Simpson's paradox shows that if we interpret the probabilities as causes we may make wrong decisions. As an example, shows the recovery rate of a drug test in which depending on the groups separation the most probable outcome changes.

1. The probability of recovery is higher for patients that took the drug

$$p(\text{recovery}|\text{drug}) > p(\text{recovery}|\text{no drug}) \quad (7.1)$$

2. For both male and female patients the relation was oposite

$$p(\text{recovery}|\text{drug, male}) < p(\text{recovery}|\text{no drug, male}) \quad (7.2)$$

$$p(\text{recovery}|\text{drug, female}) < p(\text{recovery}|\text{no drug, female}) \quad (7.3)$$

endogenous variables are binary variables that we are interested in.

Exogenous variables are latent, independent binary variables that affect externally the state of our endogeneous variables.

A *Structural Causal Model (SCM)*, also known as *Structural Equation Model (SEM)*, is a tuple ...

There is one Structural equations per endogenous variable.

1. a product of standard measures ...
2. a product of standard measures ...
3. Measurable mapping ...
4. A product probability measure. ...

An augmented functional graph $G^a(M)$ depicts the exogenous variables while the functional graph $G(M)$ doesn't.

If M has no *self-loops*, the causal graph of M is a subgraph of the functional graph $G(M)$.

Definition 7.1. We call the family of sets of probability distributions of the solutions of $M_{do(I, E_I)}$...

some of the previous points are the basic difference between causal models and probabilistic models.

We denote a marginalization of the model M with respect two variables X_2 and X_4 as $M \setminus \{2, 4\}$.

See the following extra references De Mooij (2013), and Bongers and Mooij (2018), Blom and Mooij (2018), Bongers et al., 2018.

Definition 7.2. Definitions of: Independence and conditional independence

Definition 7.3. Definition of: nodes blocking a path

Theorem 7.4. For an acyclic SCM, ...

Reichenbach's principle of common cause, the dependence $X|Y$

The Reichenbach's Principle may fail in case of *selection bias* (related with the explaining away problem)

7.3 Causal Inference: Predicting Causal Effects

Theorem 7.5. *Back-Door Criterion Pearl (2000)*

7.4 Resolving Simpson's paradox

It is important to realize that “seeing” is not the same as “doing”.

- $p(R = 1|D = 1)$: the probability that somebody recovers, given the observation that the person took the drug.
- $p(R = 1|do(D = 1))$: the probability that somebody recovers, if we force the person to take the drug.

In practice randomized control trial

A **randomized controlled trial** (or randomized control trial; [2] RCT) is a type of scientific (often medical) experiment which aims to reduce bias when testing a new treatment. The people participating in the trial are randomly allocated to either the group receiving the treatment under investigation or to a group receiving standard treatment (or placebo treatment) as the control. Randomization minimises selection bias and the different comparison groups allow the researchers to determine any effects of the treatment when compared with the no treatment (control) group, while other variables are kept constant. The RCT is often considered the gold standard for a clinical trial. RCTs are often used to test the efficacy or effectiveness of various types of medical intervention and may provide information about adverse effects, such as drug reactions. Random assignment of intervention is done after subjects have been assessed for eligibility and recruited, but before the intervention to be studied begins. – **Wikipedia**

7.5 Causal Discovery: from data to causal graph

Randomized controlled trials Fisher (1935) are one solution to avoid previously seen problems.

1. Divide patients into two groups: treatment and control randomly
2. Patients with the treatment group are forced to take a drug, and patients in the group are forced to not take the drug (but to take a placebo instead): $D = C$
3. Estimating the causal effect of the drug now becomes a standard ...
4. ...

7.5.1 Local Causal Discovery (LCD)

Simple method that Joris Mooij likes

7.6 Practical application

It is possible to apply k -fold Cross-validation to the observational data and interventional data in order to estimate the test performance.

7.7 Conclusions

Additional readings: Causality: Models reasoning and inference, pearl 2000. Constraint-based causal discovery for non-linear

- Elements of causal inference, foundations and learning algorithms by Peters, Janzing, Scholkopf 2017
- Causation, prediction and search by spirtes, glymour, scheines 2000
- correlation and causation by Wrights 1921
- Causal inference in statistics: an overview by Pearl 2009
- Simpson's paradox: an anatomy by Pearl 1999
- Causality: models, reasoning and inference by Pearl 2000
- Theoretical aspects of cyclic structural causal models by Bongers, Peters, Scholkopf, Mooij 2018
- Markov properties for graphical models with cycles and latent variables by Forré, and Mooij 2017

Some possible applications of Causal inference could be:

- Transfer learning
- Domain adaptation
- Reinforcement learning

What tools or frameworks: **there are no tools yet**, R package for individual methods. Literature and the implementations are scattered, **necessary to unify!**.

Bibliography

- Blom, T. and Mooij, J. M. (2018). Generalized structural causal models. *arXiv preprint arXiv:1805.06539*.
- Bongers, S. and Mooij, J. M. (2018). From random differential equations to structural causal models: the stochastic case. *arXiv preprint arXiv:1803.08784*.
- De Mooij, M. (2013). *Global marketing and advertising: Understanding cultural paradoxes*. Sage Publications.
- Fisher, R. A. (1935). The design of experiments.
- Pearl, J. (2000). Causal inference without counterfactuals: Comment. *Journal of the American Statistical Association*, 95(450):428–431.

Chapter 8

Reinforcement Learning by Jan Peters

Fri. 14:30–16:00

8.1 Optimal Control Systems

Data to model to value function to policy back then to data.

8.1.1 Markov Decision Problems

A stationary MDP is defined as

- a state space $s \in S$
- action space $a \in A$
- transition dynamics $P(s_{t+1}|s_t, a_t)$
- reward function $r(s, a)$
- initial state probabilities $\mu_0(s)$

The Markov property says that the transition dynamics depends only on the current time step.

$$P(s_{t+1}|s_t, a_t, s_{t-1}, a_{t-1}, \dots) = P(s_{t+1}|s_t, a_t) \quad (8.1)$$

8.1.2 Basic reinforcement learning loop

The objective is to maximize the expected long-term reward

$$J_\theta = \mathbb{E}_{\mu_0, P, \pi} \left[\sum_{t=1}^{T-1} \gamma^t r(s_t, a_t) \right] \quad (8.2)$$

The algorithmic description of the value iteration

- Init: $V_T^*(s) \leftarrow r_T(s), t = T$
 - Compute Q-Function for time step t (for each state action pair)

$$Q_t^*(s, a) = r_t(s, a) + \gamma \sum_{s'} P_t(s'|s, a) V_{t+1}^*(s') \quad (8.3)$$

- Compute V-Function for time step t (for each state) (TODO, check if next equation is max or argmax)

$$V_t^*(s) = \max_a Q_t^*(s, a) \quad (8.4)$$

- Repeat: $t = t - 1$

- Until $t = 1$
- RReturn: Optimal policy for each time step

$$\pi_t^*(s) = \arg \max_a Q_t^*(s, a) \quad (8.5)$$

The Bellman Equation (Bellman Principle of Optimality)

$$V^*(s) = \max_a (r(s, a) + \gamma \mathbb{E}_p[V^*(s')|s, a]) \quad (8.6)$$

See Policy evaluation with temporal differences: a survey and comparison Dann et al. (2014)

When the max is expensive is possible to use the *Policy Iteration* method:

1. Policy evaluation: Estimate quality of states (and actions) with current policy
2. Policy improvement: Improve policy by taking actions with the highest quality

8.1.3 Linear Quadratic Gaussian Systems

A Linear Quadratic Regulator (LQR) system is defined as

- state space *in* \mathbb{R}^n
- action space *in* \mathbb{R}^m
- linear transition dynamics with Gaussian noise

$$p_t(x_{t+1}|x_t, u_t) = \mathcal{N}(x_{t+1}|A_t x_t + B_t u_t + b_t, \Sigma_t) \quad (8.7)$$

- quadratic reward function

$$r_t(x, u) = (x - r_t)^T R_t (x - r_t) + u_t^T H_t u_t \quad (8.8)$$

$$r_t(x) = (x - r_T)^T R_T (x - r_T) \quad (8.9)$$

- initial state density

$$\mu_0(x) \mathcal{N}(x|\mu_0, \Sigma_0) \quad (8.10)$$

See Stefan Schaal, Christopher G. Atkeson 1998 Schaal and Atkeson (1998)

The LQR systems need the initial point to be linearly related to the optimal point (eg. keeping a stick balanced upwards). However, they can not solve situations in which it is necessary a non-linear component (eg. if the stick starts from a hanging position, and it needs some sinoidal movement before it can reach the upward position).

See work by Emo Todorov and Yuval Tassa on Incremental LQG (a simplification of Differential Dynamic Programming by Dyer and McReynolds 1969 Dyer and McReynolds (1969))

With Optimal control we can compute optimal policies but only on

1. Discrete Systems: (but world is not discrete)
2. Linear Systems, Quadratic Reward, Gaussian Noise (LQR): (but the world is not linear)
3. Along an optimal trajectory: (it is really hard to find)

For these reasons we need to approximate.

8.2 Value Function Methods

Data to value function to policy to data.

One of the principles is that “all models are wrong, but some are useful”. In the previous section we created perfect models, but they may not be match the real world. In that case, there can be drastical problems.

The *Classical Reinforcement Learning* postulate is to solve the optimal control problem by learning the value function and not the model.

8.2.1 Markov Decision Processes (MDP)

Infinite Horizon with a discounted reward parameter γ

Updates the value function based on samples $D = \{s_i, a_i, r_i, s'_i\}$

8.2.2 Temporal difference learning

We are incorporating an error value into the prediction of the states (see Reinforcement learning, rich sutton and andy barto, 1998 Sutton et al. (1998)).

With our estimate we can compute the TD error and make a decision

- if the estimation was higher, we decrease the prediction
- if the estimation was lower, we increase the prediction

1. Init $V_0^*(s) \leftarrow 0$
2. Repeat $t = t + 1$
 - Observe transition (s_t, a_t, r_t, s_{t+1})
 - Compute TD error $\delta_t = r_t + \gamma V_t(s_{t+1}) - V_t(s_t)$
 - Update V-Function $V_{t+1}(s_t) = V_t(s_t) + \alpha \delta_t$

3. until convergence of V

But we do not want deterministic policies as these will not explore the space, for that reason there are at least two policy for exploration

1. Epsilon-Greedy Policy
2. Soft-Max Policy (it has an important temperature parameter β)

Update equations for learning the Q-function $Q(s, a)$

$$Q_{t+1}(s_t, a_t) = \dots \quad (8.11)$$

In which it is necessary to estimate the future action a_t . There are two methods

1. Q-learning: $a_t = \arg \max_a Q_t(s_{t+1}, a)$
2. SARSA: $a_t = a_{t+1}$, where $a_{t+1} \sim \pi(a|s_{t+1})$

8.2.3 Approximating the value function

Instead of creating the matrix V we can approximate with any function approximation method (see Dann et al: Policy evaluation with temporal differences: a survey and comparison, JMLR, 2014 Dann et al. (2014)).

Some remarks on temporal difference learning

- It is not a proper stochastic gradient descent
- As the target values $V^\pi(s)$ change after each parameter update
- This “ignorance” introduces a bias in our optimization
- ...

8.2.4 Batch-Mode Reinforcement Learning

Online methods are typically data-inefficient as they use only once every sample. The reuse of the samples has been done in Least-squares temporal difference learning and fitted q-iteration (Tree-Based batch mode reinforcement learning Ernst et al. (2005), Batch reinforcement learning Lange et al. (2012)).

In Q-iteration we do as in Value-iteration, but we use Supervised Learning methods to approximate the Q function.

See Reinforcement learning in robot soccer, 2009 Riedmiller et al. (2009).

8.3 Policy Search

From data to policy and back to data.

see Reinforcement learning of motor skills with policy gradients, 2008 Peters and Schaal (2008).

8.3.1 Black-box approaches

- Perturb the parameters of your policy: $\delta J = J(\theta + \delta\theta) - J(\theta)$
- Approximate J by the first order Taylor approximation $J(\theta + \delta\theta) = J(\theta) + \frac{\partial J(\theta)}{\partial \theta} \delta\theta$
- Solve for $\frac{\partial J(\theta)}{\partial \theta}$ in a least squares sense (linear regression).

See a large class of algorithms includes: Kiefer-Wolfowitz procedure, Robbins-Monroe, Simultaneous Perturbation Stochastic Approximation SPSA,...

8.3.2 Likelihood-Ratio Policy Gradient methods

The expected long term reward $J(\theta)$ can be written as expectation over the trajectory distribution.

8.4 Key problems

1. no notion of data in the generic problem formulation
2. optimization bias problematic with data
3. role of features is unclear in most methods

Bibliography

- Dann, C., Neumann, G., and Peters, J. (2014). Policy evaluation with temporal differences: A survey and comparison. *The Journal of Machine Learning Research*, 15(1):809–883.
- Dyer, P. and McReynolds, S. (1969). Optimization of control systems with discontinuities and terminal constraints. *IEEE Transactions on automatic Control*, 14(3):223–229.
- Ernst, D., Geurts, P., and Wehenkel, L. (2005). Tree-based batch mode reinforcement learning. *Journal of Machine Learning Research*, 6(Apr):503–556.
- Lange, S., Gabel, T., and Riedmiller, M. (2012). Batch reinforcement learning. In *Reinforcement learning*, pages 45–73. Springer.
- Peters, J. and Schaal, S. (2008). Reinforcement learning of motor skills with policy gradients. *Neural networks*, 21(4):682–697.
- Riedmiller, M., Gabel, T., Hafner, R., and Lange, S. (2009). Reinforcement learning for robot soccer. *Autonomous Robots*, 27(1):55–73.
- Schaal, S. and Atkeson, C. G. (1998). Constructive incremental learning from only local information. *Neural computation*, 10(8):2047–2084.
- Sutton, R. S., Barto, A. G., et al. (1998). *Reinforcement learning: An introduction*. MIT press.

Chapter 9

Probabilistic Numerics: Nano-machine-learning by Michael A Osborne

9:30–11:00

Numerics is becoming one of the important fields of Machine Learning.

In numerical analysis, numerical integration constitutes a broad family of algorithms for calculating the numerical value of a definite integral, and by extension, the term is also sometimes used to describe the numerical solution of differential equations. This article focuses on calculation of definite integrals. The term numerical quadrature (often abbreviated to quadrature) is more or less a synonym for numerical integration, especially as applied to one-dimensional integrals. Some authors refer to numerical integration over more than one dimension as cubature;^[1] others take quadrature to include higher-dimensional integration. – Wikipedia

The answer to a numeric problem can only be approximated, e.g.

$$F = \int_{-3}^3 f(x) dx \tag{9.1}$$

for

$$f(x) = \exp(-(\sin(3x))^2 - x^2) \tag{9.2}$$

This can be as long as 30 polinomic elements. Although we can not use it analytically, we can get an approximated answer with a computer.

A motivation:

1. numeric *error* is significant
2. numeric methods are generic
3. our numerics problems tax our computation

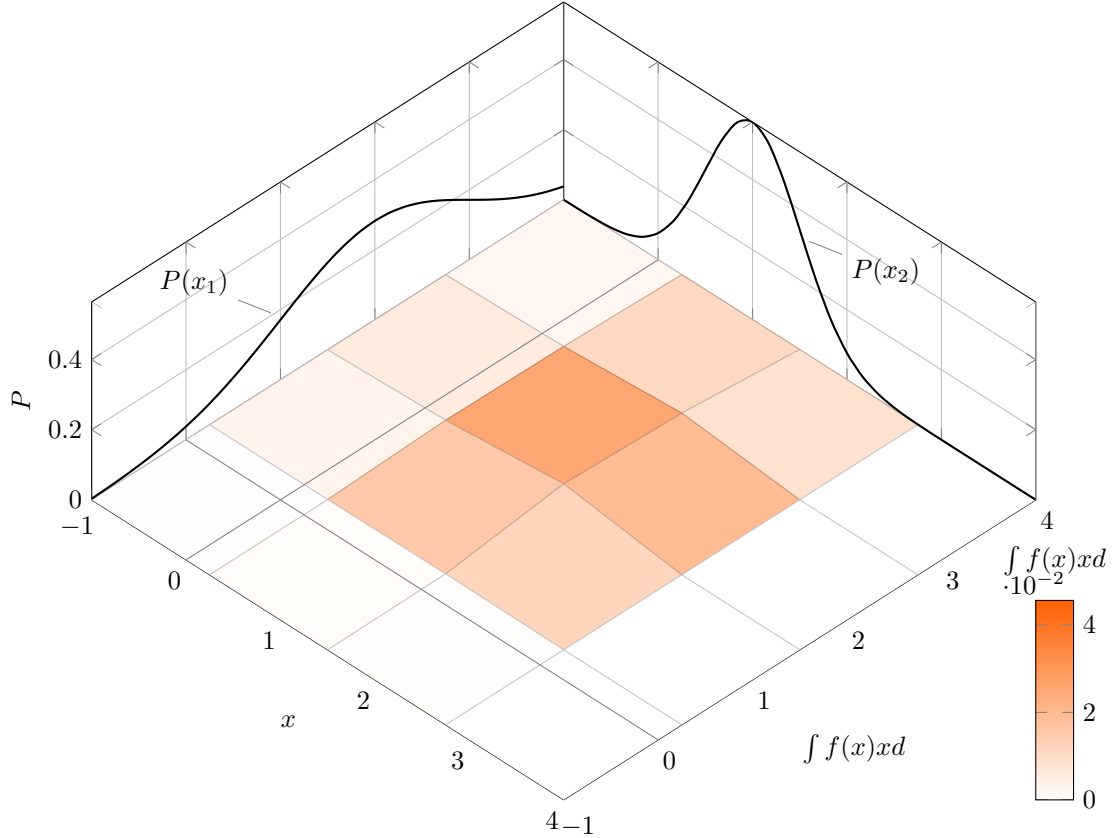
Some important parts of the integration to have in mind:

- The data: are the evaluations, or
- Predictand: is the integral
- Decisions: ...

Bayesian quadrature is probabilistic numerics for intergration

Bayesian Quadrature is a statistical approach to the numerical problem of computing integrals and falls under the field of probabilistic numerics. It can provide a full handling of the uncertainty over the solution of the integral expressed as a Gaussian Process posterior variance. It is also known to provide very fast convergence rates which can be up to exponential in the number of quadrature points n . [5] – Wikipedia

With a *Gaussian process* prior for the integrand, the *integral is joint Gaussian*.



Monte Carlo is also Bayesian quadrature
The motnte caro estimate is

$$\int f(x)p(x)dx \simeq 1/N \sum_{i=1}^N f(x_i) \quad (9.3)$$

is a maximum a-posteriori under the (imporper) prior

$$p(f) = \lim_{e \rightarrow 0} GP(0, \theta^2 \mathbb{I}(x = x') + c^{-1}) \quad (9.4)$$

TODO missing conclusion for last equation

Managing parameters θ requires the model ...

$$\int f(x|\theta)d\theta \text{ missing equation} \quad (9.5)$$

In optimization it is not enough looking for the higher likelihood value, but we are looking for the highest mass. This is picks with larga areas around.

$$p(data) \simeq \text{missing} \quad (9.6)$$

We prefer *flat optima* to *pick optima* precisely because of the mass.
In Monte Carlo estimator

$$\int f(x)p(x)dx \simeq 1/N \sum_{i=1}^N f(x_i) \quad (9.7)$$

ignores relevant information and assumes certain sample distribution (see Monte Carlo is Fundamentally unsound O'Hagan (1987))

Warped sequential active Bayesian integration (WSABI) uses a loss that is the uncertainty in the integrand (see Sampling for inference in Probabilistic models with fast bayesian quadrature, NIPS Gunter et al. (2014))

We propose a novel sampling framework for inference in probabilistic models: an active learning approach that converges more quickly (in wall-clock time) than Markov chain Monte Carlo (MCMC) benchmarks. The central challenge in probabilistic inference is numerical integration, to average over ensembles of models or unknown (hyper-)parameters (for example to compute the marginal likelihood or a partition function). MCMC has provided approaches to numerical integration that deliver state-of-the-art inference, but can suffer from sample inefficiency and poor convergence diagnostics. Bayesian quadrature techniques offer a model-based solution to such problems, but their uptake has been hindered by prohibitive computation costs. We introduce a warped model for probabilistic integrands (likelihoods) that are known to be non-negative, permitting a cheap active learning scheme to optimally select sample locations. Our algorithm is demonstrated to offer faster convergence (in seconds) relative to simple Monte Carlo and annealed importance sampling on both synthetic and real-world examples. – Abstract from Gunter et al. (2014)

In global optimization

- Data: evaluation
- Predictand: minimizer
- Decisions: location

Bayesian optimisation is probabilistic numerics for global optimization
The loss for optimisation could be

1. the lowest evaluation (value): Some times it is really important to choose the best option, but not care about the uncertainty
2. the uncertainty in the minimiser (location-information): e.g. if you create the model in synthetic data, it is important to evaluate the uncertainty on the new real data
3. the uncertainty in the minimum (value-information):

•

$$\lambda_{VL} = y_N \quad (9.8)$$

•

$$\lambda_{LIL} = \mathbb{M}(x_* | x_N, y_N, D_N) \quad (9.9)$$

Solving the intrinsic “myopia” of bayesian optimization methods (see GLASSES: relieving the myopia of bayesian optimization González et al. (2016))

9.1 Upper confidence bound

is the myopic acquisition function

$$\text{missing equation} \quad (9.10)$$

given a surrogate with mean $m(x_n)$ and variance $V(x_n)$

9.2 Information-theoretic methods

give alternative myopic implementations of va-ue.-information and location-information losses:
these methods tend to be more exploratory

9.3 Technology at work: The future of automation Tue. 9:30–11:00

What are humans still good for?

Bibliography

- González, J., Osborne, M., and Lawrence, N. (2016). Glasses: Relieving the myopia of bayesian optimisation. In *Artificial Intelligence and Statistics*, pages 790–799.
- Gunter, T., Osborne, M. A., Garnett, R., Hennig, P., and Roberts, S. J. (2014). Sampling for inference in probabilistic models with fast bayesian quadrature. In *Advances in neural information processing systems*, pages 2789–2797.
- O’Hagan, A. (1987). Monte carlo is fundamentally unsound. *The Statistician*, pages 247–249.

Chapter 10

Kernel methods by Arthur Gretton

Tue. 11:30–13:00

- Testing goodness of fit: Given a model P and samples Q
- Dependence testing

Maximum mean discrepancy (MMD)

10.1 Reproducing Hilbert spaces

Definition (Inner product)

Let H be a vector space over \mathbb{R} . A function $k(x, x')$ is an inner product of x .

Theorem 10.1. *Sums of kernels are kernels: Given $\alpha > 0$ and k, k_1 and k_2 all kernels on χ , then αk and \dots are kernels.*

Theorem 10.2. *Products of kernels are kernels:*

Theorem 10.3. *Polynomial kernels:*

A famous infinite feature space kernel, the exponentiated quadratic kernel

$$k(x, x') = \exp\left(-\frac{\|x - x'\|^2}{2\sigma^2}\right) \quad (10.1)$$

Smoothness in RKHS with exponentiated quadratic kernel

$$f(x) = \sum_{i=1}^m \alpha_i k(x_i, x) = \sum_{l=1}^{\infty} f_l[\sqrt{\lambda_l} e_l(x)] \quad (10.2)$$

10.2 Interlude: divergence measures

Integral probability metrics (subtraction): wasserstein, MMD, TV

$$D_H(P, Q) = \sup_{g \in H} |E_{X \sim P} g(X) - E_{Y \sim Q} g(Y)| \quad (10.3)$$

F-divergences: Pearson chi2, Hellinger, KL, TV

$$D_f(P, Q) = \int_{\chi} q(x) f\left(\frac{p(x)}{q(x)}\right) dx \quad (10.4)$$

Notice the intersection TV (see ...)

10.3 Two-sample testing with MMD

- Null hypothesis: H_0 when $P = Q$
should see $M\hat{M}D^2$ close to zero
- Alternative hypothesis H_1 when $P \neq Q$
should see $M\hat{M}D^2$ far from zero

Set the threshold by shuffling the data from both classes and dividing into two sets P and Q . Then computing the $M\hat{M}D^2$ of P and Q .

10.4 Training GANs with MMD Wed. 9:30–11:00

Generative Adversarial Networks are composed of a student (generator) and a teacher (discriminator). The student is learning to generate samples and the teacher is assessing if the generated samples are good or not. However, the student could memorize one sample and improve this one until the teacher assesses that it is always correct.

To improve the critic witness it is possible to add convolutional features to be discriminated, and the teacher (critic) also needs to be trained (see MMD GAN Li et al, NIPS 2017 Li et al. (2017), Coulomb GAN Unterthiner et al., ICLR 2018 Unterthiner et al. (2017))

Another idea is WGAN-GP (see Wasserstein GAN by Arjovsky et al. ICML 2017 Arjovsky et al. (2017), WGAN-GP Gulrajani et al. NIPS 2017 Gulrajani et al. (2017)).

New MMD GAN witness regulariser (Arbel, Sutherland, Binkowski, G NIPS 2018), based on semi-supervised learning regulariser by Bousquet et al NIPS 2004.

10.4.1 The kernel inception distance (KID)

The kernel inception distance (by Binkowski, sutherland, arbel G ICLR 2018) measures the similarity of the samples' representations in the inception architecture (pool3 layer) MMD with kernel

$$k(x, y) = \left(\frac{1}{d}x^T y + 1\right)^3 \quad (10.5)$$

Checks match for feature means, variances and skewness. It is unbiased?.

10.5 Testing statistical dependence

Example with captions of images of cats and dogs. First we obtain a good kernel to compare images $k(x, x')$, and a good kernel to compare text $l(x, x')$.

- Given: samples from a distribution P_{XY}
- Goal: are X and Y independent?

$$MMD^2(\hat{P}_{XY}, \hat{P}_X, \hat{P}_Y, H_k) = \frac{1}{n^2} \text{trace}(KL) \quad (10.6)$$

10.5.1 Finding covariance with smooth transformations

Illustration with a variable X and Y with a shape of a circle perimeter with some gaussian noise. In this case the correlation between variables is 0, but with certain witness functions for $w_x(X)$ and $w_y(Y)$ we can find a correlation between them.

10.5.2 Application: dependence detection across languages

- Testing task: detect dependence between English and French text
- k-spectrum kernel, $k = 10$

Bibliography

- Arjovsky, M., Chintala, S., and Bottou, L. (2017). Wasserstein generative adversarial networks. In *International Conference on Machine Learning*, pages 214–223.
- Gulrajani, I., Ahmed, F., Arjovsky, M., Dumoulin, V., and Courville, A. C. (2017). Improved training of wasserstein gans. *CoRR*, abs/1704.00028.
- Li, C.-L., Chang, W.-C., Cheng, Y., Yang, Y., and Póczos, B. (2017). Mmd gan: Towards deeper understanding of moment matching network. In *Advances in Neural Information Processing Systems*, pages 2203–2213.
- Unterthiner, T., Nessler, B., Seward, C., Klambauer, G., Heusel, M., Ramsauer, H., and Hochreiter, S. (2017). Coulomb gans: Provably optimal nash equilibria via potential fields. *arXiv preprint arXiv:1708.08819*.

Chapter 11

An introduction to Bayesian nonparametrics by Sinead Williamson

11.1 The Dirichlet process

11.1.1 An urn representation

- Conditioned on $\phi \dots$

$$p(z_i = k | z_{1:i-1}) = \int p(z_i = k | \pi) p(\pi | z_{1:i-1}) d\pi \quad (11.1)$$

$$= \frac{\sum_{j=1}^i \mathbb{I}(z_j = k) + \alpha_k}{i - 1 + \sum_k \alpha_k} \quad (11.2)$$

11.1.2 Exchangeability

- Changing the order of the observation does not change the probabilities: $p(r, g, g, r, b, r) = p(r, r, r, g, g, b)$
- This allows us to treat every data point as if it were the last one that we picked out

11.1.3 Choosing the number of clusters

- The Dirichlet distribution is a great choice when there is a clear, fixed number of clusters... but sometimes that is not the case
- The finite mixture model had K mixture components:

$$p(x_n | \pi, \{u_k\}, \{\Sigma_k\}) = \sum_{k=1}^K \pi_k \mathcal{N}(x_n | \mu_k, \sigma_k) \quad (11.3)$$

- To make sure we never run out of clusters, no matter how many data points we see, we need (countably) infinite clusters

$$p(x_n | \pi, \{u_k\}, \{\Sigma_k\}) = \sum_{k=1}^{\infty} \pi_k \mathcal{N}(x_n | \mu_k, \sigma_k) \quad (11.4)$$

11.1.4 Constructing an appropriate prior

- Start off with w elements

$$\pi^{(2)} = (\pi_1^{(2)}, \pi_2^{(2)}) \sim \text{Dirichlet}(\alpha/2, \alpha/2) \quad (11.5)$$

- Split each component according to our beta

$$\pi^{(4)} = \dots \sim \text{Dirichlet}(\alpha/4, \alpha/4, \alpha/4, \alpha/4) \quad (11.6)$$

- Keep until infinity?

$$\pi^{(K)} = \dots \sim \text{Dirichlet}(\alpha/K, \dots, \alpha/K) \quad (11.7)$$

See Ferguson 1973 Ferguson (1973)

11.2 Dirichlet process and Dirichlet marginals

11.2.1 Conjugacy of the multinomial

- We saw that dirichlet distribution was a conjugate prior of the multinomial.
- This is also true for the Dirichlet process
- Pick a partition $A_1, \dots, A_k \dots$

11.2.2 The Chinese restaurant process

- We can describe a sample from a DP-distributed probability distribution in terms of the following restaurant metaphor
- Restaurant with infinitely many tables (serving different dish)
- The first person sits in the first empty table
- Second person sits at the first table with probability $1/(1 + \alpha)$ or at a new table with $1/(1 + \alpha)$.
- let m_k be the number of people sat at the k th table. The n th customer sits at the k th table with probability $m_k/(1 + \alpha)$

11.2.3 The stick breaking construction

See Sethuraman 1994 Sethuraman (1994).

Imagine a stick of unit length representing the total probability

1. Sample a $Beta(1, \alpha)$ random variable b_k
2. Break a fraction b_k of the stick. This is the first atom.
3. Sample a random location for this atom
4. recurse on the remaining stick to get
5. Repeat from $k = 1, 2, \dots$

11.2.4 Indian Buffet Process

1. A customer enters a restaurant with an infinite number of dishes
2. ...

11.2.5 Bulding latent feature models using the IBP

The number of latent features (apple, skull, thread, hat, ...) can use a Indian Buffet Process (IBP).

- Unbounded number of latent features
- Each column of the IBP corresponds to one of an infinite number of features item Weach row of the IBP selects a finite subset of these features
- The rich-get-richer property of the IBP ensures features are shared between data points
- We must pick a likelihood model that determines what the features look like and how they are combined

In order to do inference in the IBP

$$Z \sim IBP(\alpha) \quad (11.8)$$

$$A_k \sim \mathcal{N}(0, \theta_A^2 \mathbf{I}) \quad (11.9)$$

$$x_n \sim \mathcal{N}(z_n A^T, \sigma_X^2 \mathbf{I}) \quad (11.10)$$

11.2.6 Summary

- The Dirichlet process is an infinite-dimensional analoge of the Dirichlet distribution
- we use the Dirichlet distribution for clustering data into K clusters item similarly, we can use the Dirichlet process to cluster data into an unbounded (and growing) number of clusters
- the indian buffet process is an infinite-dimensional model fro feature subset selection
- we can use it to construct latent feature models with infinitely many features
- we can customize the latent feature model to match our data
- many more building blocks –gamma process, poisson process, pitman-yor process, kingman’s coalescent
- Next, we will take a look at some hierarchical models that use the DP and IBP as building blocks

11.2.7 Latent Dirichlet allocation

Dirichlet distributions are commonly used in topic models. These models describe documents using a distribution over the “topics”, where each topic is a distribution over words (see Latent Dirichlet Allocation by Blei et al. 2003 Blei et al. (2003))

1. For each topic $k = 1, \dots, K$ sample a distribution over the words $\beta_k \sim \text{Dirichlet}(\eta_1, \dots, \eta_V)$

11.2.8 Hierarchical Dirichlet process

$$G_0 \sim DP(\gamma, H) \quad (11.11)$$

$$G_m \sim DP(\alpha, G_0) \quad (11.12)$$

With small values of α the topics will be differnt, however with large α all the models will be the same.

11.2.9 The Chinese restaurant franchise

1. First restaurant (document)
2. Customers pick tables according to a Chinese restaurant process with parameter α
3. Each table asks their waiter to pick a dish
4. The waiter considers all the dishes that have been served previously in the franchise
 - Since it is the first restaurant the first table gets a new dish
 - Second table gets the previous dish with probability $1/(1 + \gamma)$ or a new otherwise
 - Keep going like the previous example of a Chinese restaurant
5. The second restaurant
6. The costumers pick tables according to a Chinese restaurant process with parameter α (from all the possible tables seen in all the franchises)
7. ...

11.2.10 Basic network models: Erdős-Renyi models

11.3 Further resources

- A tutorial on bayesian nonparametric models, S.j. gershman and D.M. Blei
- The introduction of Erik Sudderth's PhD thesis
- Markov chain sampling methods for Dirichlet process mixture models, RM Neal, Journal of computational and graphical statistics, 2000
- Python: bnpy-dev, PyIBP
- Julia: BNP.jl

Bibliography

- Blei, D. M., Ng, A. Y., and Jordan, M. I. (2003). Latent dirichlet allocation. *Journal of machine Learning research*, 3(Jan):993–1022.
- Ferguson, T. S. (1973). A bayesian analysis of some nonparametric problems. *The annals of statistics*, pages 209–230.
- Sethuraman, J. (1994). A constructive definition of dirichlet priors. *Statistica sinica*, pages 639–650.

Chapter 12

Machine Learning and Causal Inference for (Reliable) Decision Support. by Suchi Saria

Wed. 16:30–18:00

Algorithmic fairness (eg. suitability of an applicant for a job position)

In a medical environment the label given to a patient after some of the interventions depend on the actual interventions done to the patient. Eg. if we have data from a patient before any intervention, and after the interventions we get a label about their recovery. We can not know if a similar patient is a right risk, because the label depends on the intervention.

Suchi Saria is proposing the following

- The objective is not the prediction y (final label)
- Recasting the problem as “what if”: eg. what would be the outcome of the patient if we didn’t make any intervention?

See Caruana et al., KDD 2015 Caruana et al. (2015), Schulam and Saria, NIPS 2017 Schulam and Saria (2017)
An example of the use of “what if” formalization

- You are concerned about blood pressure and if you should start to do some exercise to improve it.
- Formulatoin 1: What if I were to exercise?
- Formulation 2: What is the effect of exercise on the blood pressure of individuals like myself?
- Formulation 3: What is the effect of exercise on blood pressure?

Example: Exercise and blood pressure

Core assumptions: Positivity: Every subject has non-zero probability of receiving every treatment.

Consider the confounders : covariate that has a causal effect on both the treatment and outcome. Eg. $X \rightarrow Z \rightarrow Y \leftarrow X$

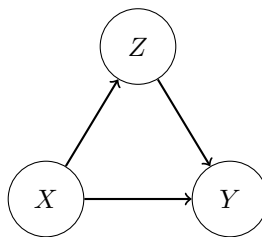


Figure 12.1: Causality graph

- Core assumptions: no unobserved confounders

One solution is the randomized control trial that removes one of the causal relations between a possible confounder X and Z . Leaving the causal connections $X \rightarrow Y \leftarrow Z$

- Randomize trials may be impossible
- In many cases we can collect observational data
- **Assumptions are not always testable from data**
- **No escape:** must rely on domain knowledge

12.1 Observed confounders

12.1.1 Feature Matching

- Search for all the matches in your data between all the features except the one being studied (pairs of treated and untreated individuals who are very similar or even identical to each other).
- Using any distance metric between sample inputs

See Sharma and Kiciman 2018, and Stuart 2010

Propensity score matching

In propensity score matching we use a supervised model $e(x)$ that predicts Z given X and creates the new causal graph $X \rightarrow e(x) \rightarrow Z \rightarrow Y \leftarrow X$. As $e(x)$ is a d-separator, by knowing it we make independent the X from Z .

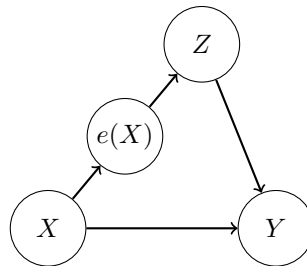


Figure 12.2: Causality graph

1. Estimate $e(x)$ using supervised learning
 - Logistic regression, or other models
 - The score must be well-calibrated
2. Distance is the difference between the propensity scores

$$Distance(x_i, x_j) = |\hat{e}(x_i) - \hat{e}(x_j)| \quad (12.1)$$

3. If the model is perfect and always predicts 1 or 0, it is not possible to match the people, as everybody will fall into the same bucket (think about assessing calibration when the model only makes predictions 0 or 1)

Weighting

12.2 References

- See Reliable decision support using conterfactual models Schulam and Saria (2017)

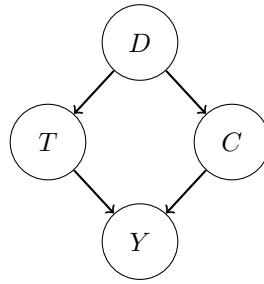


Figure 12.3: Causality graph

12.3 Unstable paths

Learning $T|C, Y$ is unstable because C still depends on an unknown variable D .

Consider naive discriminative model $P(T|C, Y)$, in this case C is **vulnerable** because it has an active unstable path to T .

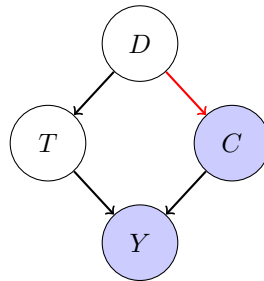


Figure 12.4: Vulnerable variables

We can create a new feature tha ...

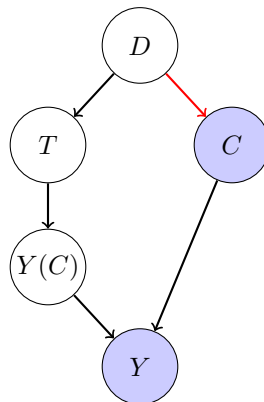


Figure 12.5: Vulnerable variables

Bibliography

- Caruana, R., Lou, Y., Gehrke, J., Koch, P., Sturm, M., and Elhadad, N. (2015). Intelligible models for health-care: Predicting pneumonia risk and hospital 30-day readmission. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 1721–1730. ACM.
- Schulam, P. and Saria, S. (2017). Reliable decision support using counterfactual models. In *Advances in Neural Information Processing Systems*, pages 1697–1708.

Chapter 13

Machine Learning in the industry

13.1 Real-world ML Challenges at the Scale of Banking by BBVA Data and Analytics

Up to this moment, ML in banking has focused on fraud detection.

Working at the moment on an app for expenses prediction. Hundreds of millions of time series.

Classification of transactions, they tested Word2vec with Vector of Locally Aggregated Descriptors (VLAD) pooling.

Expense forecasting with LSTM layers and a dense layer on top. Now the expenses forecasting is including uncertainty on the predictions (See Uncertainty Modelling in Deep Networks: Forecasting Short and Noisy Series Brando et al. (2018))

How to compare customers (see cleint2vec: Towards systematic baselines for banking applications Baldassini and Serrano (2018))

Recommender systems, modelling users for target advertising.

In order to create a loss functions that considers (See, A Missing Information Loss function for implicit feedback datasets Arévalo et al. (2018))

Reinforcement Learning for Fair Dynamic Pricing Maestre et al. (2018)

What is not in the books of Machine Learning when applying methods in the industry:

- Fairness
- Privacy
- Ethics
- Business
- Data acquisition/labelling
- UX/UI
- Design

13.2 Data Efficient Reinforcement Learning by PROWLER.io

Based in three teams: Probabilistic team, reinforcement learning team, and multiagent team.

Curren tReinforcement learning technology, successes using deep Q-Networks.

13.3 Lynx: real-time accurate fraud detection over massive data. Instituto de Ingenieria del conocimiento by Álvaro Barbero Jiménez

There are 8 countries using Lynx, and over 30,000 million transactions processed per year.

Their system need to make decision in a few milliseconds.

Device tries to make an operation, this is sent to the institution (bank), the institution sends a copy of the operation to Lynx, that will estimate the probability of fault in some miliseconds and send the answer back to the institution. Then the institution has to choose on what action to take and send back the operation success (or deny) to the device.

The decisions are taken with two types of analysis

1. Parametric Analysis
2. ...

Approach for programming: Use Python, R, elastic and cocker in order to test ideas fast, and if the results are good, implement it on C, Bash (see book Data science at the command line) or Fortran.

One of the problems of detecting fraud is that the fraudulent transactions evolve as they are detected. It is necessary to use Adaptive models and train with incremental learning.

Hardware specification: 384 threads, 6TB RAM, 40 TB SSD (a training in one day) Training data 800k-5M transactions.

As an error measure they use Value Detection Rate.

13.4 Microsoft Research by Sebastian Nowozin

120 researchers at Cambridge (200 total workers). Cambridge focus on Machine Learning for healthcare. New lab in Montreal focused on Reinforcement learning and Deep Learning.

Why is Artificial Intelligence growing? (1) Massive computation power (GPUs, FPGAs, ...), (2) Powerful algorithms (2012 Alexnet?), (3) Big data.

AI principles

- Fairness: we want the algorithms to avoid systematic biases. It is difficult to remove biases that are already in the datasets.
- Accountability:
- Transparency:
- Ethics:

If you know what you are doing, then you are not doing research. (not sure if this was the exact wording) – Sebastian Nowozin

13.4.1 Timelines

In images, videos and audio: 2000 basic research in audio, skeletal tracking and facial recognition, 2010 Kinect, 2011 Kinect Fusion, 2012 HoloDesk, 2015 HoloLens.

In translation: 1991 basic research in natural language and speech recognition, 2007 Big bets with product team, 2014 promise of speech recognition with translation, 2015 Skype translator launches, 2016 Microsoft translator API and personal universal translator launches.

Visual Studio IntelliCode (AI-assisted development)

InnerEye, how can computers understand the segments of medical images. Previously, a doctor had to spend around 8 hours segmenting a tumour in several images. Now, a program can do the segmentation really fast, and the doctor can check the results much faster.

Adaptive, learning to decode the immune system to diagnose disease. From a blood sample, immunosequencing extracts some features from t-cells?, a Machine learning model can be trained in order to improve the health care service.

Bibliography

- Arévalo, J., Duque, J. R., and Creatura, M. (2018). A missing information loss function for implicit feedback datasets. *arXiv preprint arXiv:1805.00121*.
- Baldassini, L. and Serrano, J. A. R. (2018). client2vec: Towards systematic baselines for banking applications. *arXiv preprint arXiv:1802.04198*.
- Brando, A., Rodríguez-Serrano, J. A., Ciprian, M., Maestre, R., and Vitrià, J. (2018). Uncertainty modelling in deep networks: Forecasting short and noisy series. *arXiv preprint arXiv:1807.09011*.
- Maestre, R., Duque, J., Rubio, A., and Arévalo, J. (2018). Reinforcement learning for fair dynamic pricing. *arXiv preprint arXiv:1803.09967*.

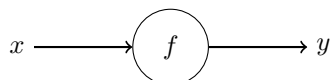
Chapter 14

Generative Adversarial Networks by Sebastian Nowozin

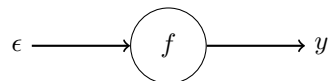
Fri. Sep 07, 9:30–11:00

14.1 Probabilistic models

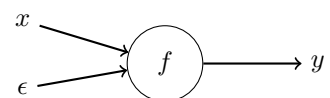
Non-probabilistic:



Probabilistic Generative:



Probabilistic discriminative (Conditionally Generative):



14.2 Example applications of GANs

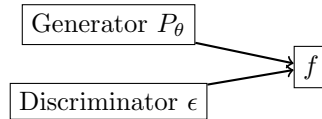
TODO: check the following graph

Examples of the evolution on the generation of human faces from 2014 to 2018: Goodfellow et al 2014 Goodfellow et al. (2014), Radford et al. 2015 Radford et al. (2015), Roth et al, 2017, Karras et al., 2018 Karras et al. (2017).

DCGAN architecture to generate the interior of rooms. This architecture is a 100 hidden representation z ? that is projected and reshaped into a $4x4x1024$ layer, then a convolution 1 of $8x8x512$, then stride 2 and kernel $5x5$ and convolution 2 of $16x16x256$ then stride 2 (See more about linear interpolation in the latent space in Radford et al., 2015 Radford et al. (2015)).

Another example from the same publication Radford et al. (2015) of Vector Arithmetics in the hidden vector space. In the example the authors compute the mean image of man with glasses, then subtract the mean representation of man without glasses and sum the mean representation of woman without glasses, with the result of new generated images of woman with glasses.

Another example is the Image Super-Resolution by Ledig et al., CVPR 2017 Ledig et al. (2017) in which the authors show a method to give as an input a low resolution image to a GAN and it augments the resolution of the image (outperforming previous state-of-the-art methods).



14.3 Principles of estimation

The classic parametric models (eg. fitting a Gaussian) use a density function, have a limited expressive power in a limited number of parameters (eg. mean and variance), and it is a mature field ...

One of the best known methods is the Likelihood and Maximum Likelihood Estimation (MLE) Fisher 1929 Fisher (1929)?.

An example to maximize the likelihood of data given that the model is a Gaussian. If we assume that every point is independent, we want the probability of all the points being maximized.

$$L(\theta) = \prod_i p(x_i|\theta) \quad (14.1)$$

$$\log L(\theta) = \log \prod_i p(x_i|\theta) \quad (14.2)$$

$$\log L(\theta) = \sum_i \log p(x_i|\theta) \quad (14.3)$$

14.3.1 Implicit models

Three important publications from 1990 till now ...

- Problem 1: Non-invertible Map. Nothing guarantees that the generative mapping is invertible.
- Problem 2: Lack of Density. We are mapping a low dimensional space μ to a high dimensional space p with a continuous function. This means that in the output space p all the density is in a small manifold, thus having points where $p(x)$ is not defined a.e.
- Problem 3: Misspecification: (see White 1994, Estimation, Inference, and specification analysis White (1996). The following is an elementary example:
 - We know that a coin is biased.
 - Prior: we have a uniform $p(a) = p(b) = 1/2$ (heads and tails)
 - Shows that the posterior under a fair coin with number of draws $k \in \{1, \dots, 60\}$ keeps oscillating towards a and b occasionally being in plateaus.

14.4 GAN models

A list of all the named GANs The GAN Zoo

A division of GANs space

- Not defined in the dimensionally misspecified case: KL, f-GAN, JS-GAN
- Defined in the dimensionally misspecified case: Generalized f-GAN, IPMs, Wasserstein MMD, μ -fisher...

14.4.1 Divergences and f-GAN family

We want to optimize Saddle Point objectives.

- Practical difficulty: non monotonic objective
- Theoretical difficulties: which algorithm to use?

See more in Nowozin, Cseke, Tomioka NIPS f-GAN 2014 Nowozin et al. (2016)
Explanation of f -divergences

- Divergence between two probability densities (See Csiszar and Shields, 2004, and Liese and Vajda, IEEE Inf Th, 2006)
- Scalar “generation function”
- Assumptions: Both distributions have a density function wrt Lebesgue measure

See how to estimate f -divergences from samples in Nguyen, Wainwright, Jordan, 2010 Nguyen et al. (2010).
Here is the resulting lowerbound

$$D_f(P||Q) \geq \sup_{T \in \mathcal{T}} (\mathbb{E}_{X \sim P} [T(x)] - \mathbb{E}_{X \sim Q} [f^*(T(x))]) \quad (14.4)$$

Where the first expectation is approximated using samples from P and the second with samples from Q .
Then we can compare the objective of a GAN

$$\min_{\theta} \max_w (\mathbb{E}_{X \sim P_{\theta}} [\log D_w(x)] + \mathbb{E}_{X \sim Q} [\log(1 - D_w(x))]) \quad (14.5)$$

with the objective of the f -GAN

$$\min_{\theta} \max_w (\mathbb{E}_{X \sim P_{\theta}} [T_w(x)] - \mathbb{E}_{X \sim Q} [f^*(T_w(x))]) \quad (14.6)$$

Key properties

- Invariance to coordinate transformations
- Come in pairs:

$$D_f(P||Q) = D_g(Q||P) \quad (14.7)$$

$$g(u) = uf(1/u) \quad (14.8)$$

14.4.2 IPM

- F class of real-valued bounded measurable functions
- P, Q probability measures

$$Y_F(P, Q) = \sup_{f \in F} \left| \int f dP - \int f dQ \right| \quad (14.9)$$

Choice of F determines the metric

- See more in Sriperumbudur et al., 2009 and Sriperumbudur et al., 2012.

14.4.3 IPM family: MMD

Reproducing Kernel Hilbert Space (RKHS) Norm

See Gretton et al., “A kernel two-sample test” JMLR 2012 Gretton et al. (2012)

$$Y_F(P, Q) = \sup_{f \in F} \left| \int f dP - \int f dQ \right| = \|\mu_P - \mu_Q\|_H \quad (14.10)$$

Kernel MMD training in deep learning

See more in

- Deep generative models Li et al., 2015 Li et al. (2015), and Dziugaite et al., 2015 Dziugaite et al. (2015).
- Deep discriminative models “Disco net” Bouchacourt et al., NIPS 2016 Bouchacourt et al. (2016)
- use for model criticism Sutherland et al., ICLR 2017 Sutherland et al. (2017)
- more discriminative kernel functions Li et al., NIPS 2017 Li et al. (2017)

14.4.4 IPM family: Wasserstein GANs

In computer vision look for earth movers distance

$$W(P, Q) = \inf_{U \in \Pi(P, Q)} \mathbb{E}_{(x, y) \sim U} [\|x - y\|] \quad (14.11)$$

Kantorovich-Rubinstein Duality, the previous equation is the same as

$$W(P, Q) = \sup_{\|f\|_{L \leq 1}} \mathbb{E}_{X \sim P} [f(x)] - \mathbb{E}_{X \sim Q} [f(x)] \quad (14.12)$$

See more in Arjovsky et al., WGAN. In which instead of having the constrain $\|f\|_{L \leq 1}$ it is constrained by a constant $\|f\|_{L \leq k}$. It guarantees the K -Lipschitz bounded functions.

However, it required the choice of a clipping value (eg. $[-0.01, 0.01]$), and leads to non-uniform bounding of gradients

, and Gulrajani et al., NIPS 2017 WGAN-GP. In which they approximate the Lipschitz condition with a soft-penalty

$$\tilde{W}(P, Q) = \mathbb{E}_{X \sim P} [f(x)] - \mathbb{E}_{X \sim Q} [f(x)] + \lambda \mathbb{E}_{X \sim M(P, Q)} [(\|\nabla_X f(x)\|_2 - 1)^2] \quad (14.13)$$

14.5 Problems and Fixes: Mode Collapse, Instability

Empirically observed behaviour where model produces only a few distinct samples. In order to solve the GAN model collapse and stability issues

From a

- Divergence viewpoint: Arjovsky and Bottou, 2016, Sonderby et al., 2016, Roth et al., 2017, Mescheder et al., 2018.
- Algorithmic viewpoint: Mescheder et al, [2017, 2018]

See Unstable training behaviour by Roth et al., 2017 in which they add a regularization

$$\min_{\theta} \max_w (\mathbb{E}_{X \sim P_{\theta}} [T_w(x)] - \mathbb{E}_{X \sim Q} [f^*(T_w(x))]) \quad (14.14)$$

$$- \frac{\gamma}{2} \mathbb{E}_{X \sim Q} [f^{**}(T_w(X)) \|\nabla T_w(x)\|^2] \quad (14.15)$$

TODO: solve problem with f^{**}

Simple gradient penalties

$$\min_{\theta} \max_w (\mathbb{E}_{X \sim P_{\theta}} [T_w(x)] - \mathbb{E}_{X \sim Q} [f^*(T_w(x))]) \quad (14.16)$$

$$- \frac{\gamma}{2} \mathbb{E}_{X \sim Q} [\|\nabla T_w(x)\|^2] \quad (14.17)$$

$$- \frac{\gamma}{2} \mathbb{E}_{X \sim P_{\theta}} [\|\nabla T_w(x)\|^2] \quad (14.18)$$

14.5.1 Spectral Normalization

See Miyato et al., ICLR 2018 Miyato et al. (2018)

Main idea is

- Limit discriminator function class to functions with bounded Lipschitz norm
- Bound global Lipschitz norm by product of Lipschitz norm per layer
- Compute Lipschitz norm per layer efficiently using power method

14.6 Implicit models more generally

14.7 Open research problems

- Quantitative Evaluation Metrics
 - Tournament as evaluation? Roth et al., NIPS 2017 Roth et al. (2017)
- GANS for discrete data
- Estimation Uncertainty
 - GANs do not have a likelihood nor a well-defined posterior
 - Early attempts, “Bayesian GAN” by Saatchi and Wilso, NIPS 2017 Saatci and Wilson (2017)
- Practical bounds on $\|f\|_L$
- New Divergences
 - μ -Fisher IPM Mroueh and Sercu, 2017
 - μ -Sobolev IPM Mroueh et al., 2017 Mroueh et al. (2017)
- Theory about Generalization
 - Generalization bounds
 - empirical study using “birthday paradox test”
 - Study of neural network distance (generalization) versus study of divergences

Bibliography

- Bouchacourt, D., Mudigonda, P. K., and Nowozin, S. (2016). Disco nets: Dissimilarity coefficients networks. In *Advances in Neural Information Processing Systems*, pages 352–360.
- Dziugaite, G. K., Roy, D. M., and Ghahramani, Z. (2015). Training generative neural networks via maximum mean discrepancy optimization. *arXiv preprint arXiv:1505.03906*.
- Fisher, R. A. (1929). Tests of significance in harmonic analysis. *Proc. R. Soc. Lond. A*, 125(796):54–59.
- Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. (2014). Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680.
- Gretton, A., Borgwardt, K. M., Rasch, M. J., Schölkopf, B., and Smola, A. (2012). A kernel two-sample test. *Journal of Machine Learning Research*, 13(Mar):723–773.
- Karras, T., Aila, T., Laine, S., and Lehtinen, J. (2017). Progressive growing of gans for improved quality, stability, and variation. *arXiv preprint arXiv:1710.10196*.
- Ledig, C., Theis, L., Huszár, F., Caballero, J., Cunningham, A., Acosta, A., Aitken, A. P., Tejani, A., Totz, J., Wang, Z., et al. (2017). Photo-realistic single image super-resolution using a generative adversarial network. In *CVPR*, volume 2, page 4.
- Li, C.-L., Chang, W.-C., Cheng, Y., Yang, Y., and Póczos, B. (2017). Mmd gan: Towards deeper understanding of moment matching network. In *Advances in Neural Information Processing Systems*, pages 2203–2213.
- Li, Y., Swersky, K., and Zemel, R. (2015). Generative moment matching networks. In *International Conference on Machine Learning*, pages 1718–1727.
- Miyato, T., Kataoka, T., Koyama, M., and Yoshida, Y. (2018). Spectral normalization for generative adversarial networks. *arXiv preprint arXiv:1802.05957*.
- Mroueh, Y., Li, C.-L., Sercu, T., Raj, A., and Cheng, Y. (2017). Sobolev gan. *arXiv preprint arXiv:1711.04894*.

- Nguyen, X., Wainwright, M. J., and Jordan, M. I. (2010). Estimating divergence functionals and the likelihood ratio by convex risk minimization. *IEEE Transactions on Information Theory*, 56(11):5847–5861.
- Nowozin, S., Cseke, B., and Tomioka, R. (2016). f-gan: Training generative neural samplers using variational divergence minimization. In *Advances in Neural Information Processing Systems*, pages 271–279.
- Radford, A., Metz, L., and Chintala, S. (2015). Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*.
- Roth, K., Lucchi, A., Nowozin, S., and Hofmann, T. (2017). Stabilizing training of generative adversarial networks through regularization. In *Advances in Neural Information Processing Systems*, pages 2018–2028.
- Saatci, Y. and Wilson, A. (2017). Bayesian gans. In *Advances in Neural Information Processing Systems*, pages 3624–3633.
- Sutherland, D. J., Tung, H.-Y. F., and Strathmann, H. (2017). Soumyajit de, aaditya ramdas, alexander j. smola, and arthur gretton. generative models and model criticism via optimized maximum mean discrepancy.
- White, H. (1996). *Estimation, inference and specification analysis*. Number 22. Cambridge university press.

Chapter 15

Advances in Machine Learning for Molecules by José Miguel Hernández-Lobato

Thu. 11:30–13:00

New molecules and materials can potentially solve important existing challenges like drug and medicine design for health care, energy production and storage, and greenhouse gas conversion., energy production and storage, and greenhouse gas conversion.

However, progress in drug and material discovery has been slow because of the cost of collecting data and making decision based on data, which require a lot of human intervention.

Currently there are plenty of available datasets with the properties of real and virtual molecules. It is also possible to simulate new molecules by estimating their properties with *density functional theory* (DFT) before they are made in the laboratory.

Some example projects are ...

“Robot scientist” speeds up drug discovery

Automated AI lab that learns and formulates hypotheses has identified promising anti-cancer and anti-malarial compounds

An artificial intelligence system – or ‘robot scientist’ – capable of screening potential drugs almost completely independently could speed up drug development, say the UK researchers who created it. The approach has already identified some promising leads, including an anti-cancer compound which also shows anti-malarial properties.

The robot scientist – named ‘Eve’ – is actually a collection of machines including several computers hooked up to the kind of automated instruments already found in many labs. ‘The idea is to automate scientific research,’ says lead author Ross King from the University of Manchester. ‘You tell the system about the area of research you’re interested in ... and then the computer has an automated way of forming novel hypotheses about that area of science. It can then design experiments to test these hypotheses and the lab robots go ahead and actually do the experiments.’ The computer can also interpret the results, modify its hypotheses and construct new tests completely autonomously, only needing occasional human assistance to top up reagents and remove waste.

Eve’s predecessor – Adam – carried out genetic experiments in yeast and became the first robotic system to independently make a scientific discovery. ‘What we wanted to do with Eve was apply the same approach to something with more immediate societal values,’ says King. Eve was designed specifically for drug development and the team initially chose to focus on neglected tropical diseases. Many parts of the process, such as high throughput library screening, are already highly automated, and Eve is kitted out with all the necessary equipment to screen tens of thousands of compounds a day and identify leads. But the system can also intelligently respond to results and create standardised assays using synthetic biology. Once it has identified a lead compound, Eve can engineer specific strains of yeast needed to screen for activity against a particular disease. This makes it both quicker and cheaper than standard drug screening methods, even those that already use automated equipment. The system already has proven successes, highlighting several drug candidates that could be ‘repurposed’. For example, it showed the anti-cancer compound TNP-470 can also attack the malarial parasite *Plasmodium vivax* by inhibiting an essential enzyme. This compound is now being looked at in Brazil, where this form of malaria is prevalent.

‘It was surprising to me how easy it was to find interesting compounds for these diseases,’ says King. ‘When I started this I assumed that the main focus of the research would be proof-of-principle.’

Andrey Rzhetsky, a geneticist and computational scientist at the University of Chicago, US, who was not involved in the study, praised the group’s work. ‘I am definitely a believer in this direction of AI [artificial intelligence] work,’ he comments.

In future, King says, robot scientists like Eve could be used by pharmaceutical companies to streamline the drug development process, or explore potential new functions for existing drugs. ‘We found some really interesting compounds,’ he says, ‘And there are even more exciting results that we have yet to report.’

References: K Williams et al, J. R. Soc., Interface, 2015, DOI: 10.1098/rsif.2014.1289 **Source:** chemistryworld by Emma Stoye

Machine Learning (ML) can accelerate and automate the discovery process. If DFT can be slow, it is possible to estimate the results with ML.

Some of the challenges is the fact that common ML methods only accept vectorial input data, and molecules are commonly represented by graphs.

Some of the representations of molecules in machine learning are: (1) molecular fingerprints, (2) SMILES, or (3) Graph neural networks (GNNs).

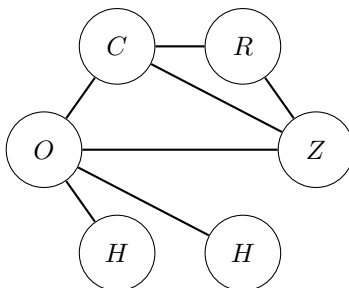
While the reaction prediction model can be achieved by Graph Neural Network (GNN)-based, or seq2seq methods.

15.1 Molecular representation for ML

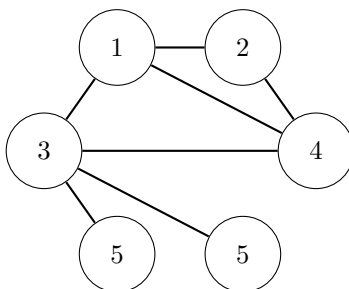
15.1.1 Molecular fingerprints

An encoding of a molecular graph into a binary string (See Rogers et al. 2010 Rogers and Hahn (2010))

The algorithm takes as an input a molecular graph with radius parameter R and length L .

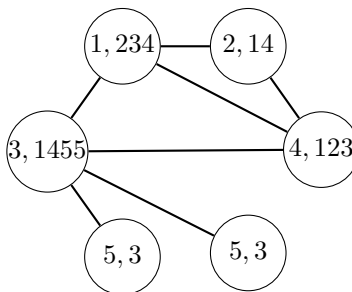


1. Assign integers to atoms by applying has function to atom features



2. For $r = 1, \dots, R$

- (a) Concatenate atom integers with integers of tneighboring atoms



- (b) Assign new integers to atoms by applying has function to concatenation

3. Create L -dimensional zero vector f
4. map generated integers to an entry in f which is set to 1. (01110001010001001000100101)

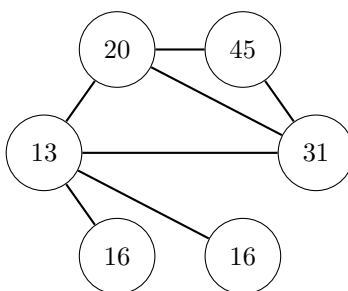
This method is fast to compute, in practice produces very good predictive performance, and it is easy to interpret as the features represent the presence of substructure.

On the other hand, the generated features are handcrafted and not data dependent, and they are not smooth (similar segments may have different representations).

15.1.2 SMILES

Simplified Molecular Input Line Entry System (SMILES) allows the representation of a molecular graph in line notation. For example, CC represents CH3CH3 (ethane), CC(=O)O represents CH3COOH (acetic acid), C1CCCCC1 represents C6H12 (cyclo...). See Sanchez-Lengleling and Aspuru-Guzik, 2018 Gómez-Bombarelli et al. (2018).

With the new sequence representations as stirngs it is possible to apply Artificial Neural Networks. Some examples are Recurrent Neural Networks applied to the string sequence to predict the following character, or use 1D convolutional Neural Networks with lefat and right context.



With this method is easy to encode molecules as simple text strings, relatively easy to understand by humans, and Natural Language Processing (NLP) methods can be applied.

One of the disadvantages of SMILES is that the representation is not invariant to the starting atom, and atoms close in the graph may be far away in the string representation.

15.1.3 Graph Neural Networks (GNNs)

This representation naturally encodes invariances to permutation of nodes, and distances between atoms (See more about GNNs in Scarselli et al. (2009), Li et al. (2018)).

A GNN includes the vectorial variables

1. $\{\vec{e}_{j \sim k}\}$ vector of edges between nodes j and k
2. $\{\vec{v}_i\}_{i=1}^N$ node features
3. \vec{u} are global features summarizing the graph properties

Set functions and auxiliary variables

Set functions have as input sets of elements and as output a single element summarizing the input set. It is invariant to input permutation and accepts a variable number of arguments. Some examples are elementwise summation, mean, maximum, minimum, ...

The forward pass in a GNN is

Given an initial set of vertices v , edges e and global features u .

1. For every edge
 - (a) updates the edge features
2. For every vertex
 - (a) summarize the incoming edge to i
 - (b) Update feature for node i
3. Summarize all edges
4. Summarize all nodes
5. Update global features
6. Compute prediction from features \vec{u}
7. return $MLP(\vec{u})$
8. Repeat all the steps several times $l = 1, \dots, L$ (every step increases the range of molecules; it is necessary to set L to capture the biggest structure in which we are interested).

Specific implementations of GNNs

- Message passing neural network (MPNN) Gilmer et al, 2017

- Gated graph neural network (GGNN) Li et al. 2016
- Weisfeiler-Lehman Network (WLN). Jin et al. 2017
- Neural graph fingerprints (NGFs). Duvenaud et al. 2015

Advantages of GNNs

Bibliography

- Gómez-Bombarelli, R., Wei, J. N., Duvenaud, D., Hernández-Lobato, J. M., Sánchez-Lengeling, B., Sheberla, D., Aguilera-Iparraguirre, J., Hirzel, T. D., Adams, R. P., and Aspuru-Guzik, A. (2018). Automatic chemical design using a data-driven continuous representation of molecules. *ACS central science*, 4(2):268–276.
- Li, Y., Vinyals, O., Dyer, C., Pascanu, R., and Battaglia, P. (2018). Learning deep generative models of graphs. *arXiv preprint arXiv:1803.03324*.
- Rogers, D. and Hahn, M. (2010). Extended-connectivity fingerprints. *Journal of chemical information and modeling*, 50(5):742–754.
- Scarselli, F., Gori, M., Tsoi, A. C., Hagenbuchner, M., and Monfardini, G. (2009). The graph neural network model. *IEEE Transactions on Neural Networks*, 20(1):61–80.

Glossary

ANN An Artificial Neural Network (ANN) is a mathematical representation of a biological neural network that simplifies its architecture and physical behaviour. It is used in Machine Learning to solve regression and decision problems. 12, 46, 48, 50, 51

GNN TODO: description. 12, 89–91