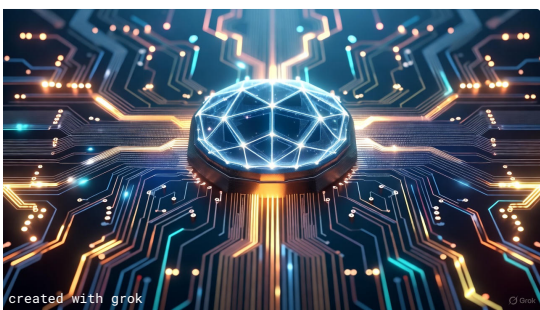


Advanced Machine Learning

Silvan Stadelmann - 6. Dezember 2025 - v0.0.1

github.com/silvasta/summary-aml



Contents

Representation	4
1 Learning objectives	4
2 Expected risk	4
3 Empirical risk	4
4 Empirical test error and expected risk	4
5 Comparing algorithm performance on test data	4
6 Data	4
6.1 Feature space	4
6.2 Example of Data	5
7 Mathematical Spaces	5
Regression	5
8 Linear Regression	5
9 Gauss Markov Theorem	5
10 Bias/Variance Dilemma	5
11 Bayesian Maximum A Posteriori (MAP) estimates	5
11.1 Ridge Regression	5
11.2 LASSO	5
11.3 Ridge vs. LASSO Estimation	5
12 Remarks on Shrinkage Methods	5
13 Learning Objectives	5
14 Gaussian Processes	6
14.1 Bayesian linear regression	6
14.2 Moments of Bayesian linear regression	6
15 Gaussian processes	6
15.1 Recall	6
15.2 Gram matrix	6
15.3 Examples of kernel functions	6

15.4	Kernel engineering by composition	6
15.5	Prediction by Gaussian processes	6
15.6	Prediction by Gaussian processes	6
15.7	Kernel validation	6
16	Controller Optimization for Robust Control	6
16.1	Gaussian processes for Control	7
16.2	Safe optimization	7
17	Model averaging is common practice	7
18	Combining Regressors - Bias	7
19	Combining Regressors - Variance	7
20	Ensemble Learning	7
21	Induction Principles for Classifier Selection	7
22	Motivation for Ensemble Methods	7
23	Weak Learners Used for Bagging or Boosting	7
24	Loss functions for classification	8
	Support Vector Machines	8
25	SVM and convex optimization	8
25.1	Learning objectives	8
25.2	Instability of the perceptron	8
25.3	Maximum-margin criterion	8
25.4	Lagrange Multipliers	8
25.5	Convex optimization	8
25.5.1	The dual	8
25.5.2	Weak duality and strong duality	8
25.5.3	Slater's condition	8
25.5.4	Consequences of strong duality	8
25.5.5	Convex optimization via the dual	8
25.6	SVMs	8
25.7	Soft-margin SVMs	8
26	Kernelization	8
26.1	Kernels	9
27	SVM summary	9
28	Structural SVMs	9
	Neural Networks	9
	Transformer	9
29	Tokenization and embeddings	9
30	Self-attention	9
30.1	Multi-headed attention	10
30.2	Cross-attention	10
30.3	Masked self-attention	10
31	Positional encodings	10
32	Broadcasting	10
33	Residual Blocks	10

34 Transformer Architecture	10
Computer Vision	10
35 CNNs and UNets	10
35.1 Convolutions and deconvolutions	10
35.1.1 Convolutional filters	10
35.1.2 Deconvolutional filters	10
35.1.3 Convolutional and deconvolutional layers	10
35.2 U-net	10
36 Stable diffusion	11
36.1 Diffusion models with images	11
36.1.1 Forward process	11
36.1.2 Reverse process	11
36.2 Conditional diffusion processes	11
37 Neural networks and variational autoencoder	11
37.1 Clip	11
38 Stable diffusion	11
38.1 Semantically relatable text embeddings with CLIP	11
38.2 Cross-attention	11
39 ViTs and MAEs	11
39.1 Vision transformers	11
39.2 Masked auto encoders	11
39.2.1 Image encoding vs text encoding	11
Diffusion	11
39.3 sgrok	12
39.4 Encoder	12
39.4.1 Clip	13
39.5 Pipeline	13
39.5.1 details	13
39.6 Unet Architecture	13
39.6.1 Multi-headed X Attention Mechanism . .	13
Graph Neural Networks	13
40 Graph	13
40.1 Neural Networks for graphs	13
40.1.1 Analyzing how neural networks work . . .	13
40.1.2 Generalization to graphs	14
40.2 Oversmoothing	14
Anomaly Detection	14
41 Data, measurements, and representations	14
41.1 Four Paradigms in Data Science	14
41.1.1 Frequentism	14
41.1.2 Bayesianism	14
41.1.3 Statistical Learning	15
41.1.4 Non-Parametric Statistics	15
41.2 Statistical learning	15
41.3 Objects, data, and representations	15
41.4 Recap	15
41.4.1 Calculus	15
41.4.2 Probability theory	15
41.4.3 information theory	15
42 Anomaly Detection	15
42.1 Dimensionality reduction with PCA	16
42.1.1 Strategy	16

42.2	Gaussian mixture models	16
42.2.1	Fitting GMMs with maximum likelihood	16
42.3	Anomaly detection with PCA and EM	16
42.3.1	Validation metric	16
42.3.2	Experimental validation	16
Reinforcement Learning		17
43	(PO-)MDP by construction	17
43.1	Policies	17
43.2	Exploration–exploitation	17
43.3	Trajectories	17
43.4	Markov property	17
43.5	Value functions	17
43.6	Bellman equations	18
43.7	Q-learning (tabular, off-policy)	18
43.8	DQN (function approximation)	18
44	Active learning	18
44.1	Information-based transductive learning	18
44.2	Safe Bayesian optimization	18
44.2.1	Bayesian optimization with ITL	18
44.3	Batch active learning	18
44.3.1	Problem formalization	18

Representation

1 Learning objectives

Estimation of Dependences Based on Empirical Data

What is the learning problem?

$$y = f_{\theta}(x) + \eta \quad \text{with} \quad \nu \sim \text{P}(\eta|0, \sigma^2)$$

2 Expected risk

- Conditional expected risk
- Total expected risk

3 Empirical risk

- Test and Train Data

Test data cannot be used before the final estimator has been selected!

Training error $\hat{R}(\hat{f}_n, \mathcal{Z}^{\text{train}})$ for Empirical Risk Minimizer (ERM)

4 Empirical test error and expected risk

Distinguish

5 Comparing algorithm performance on test data

6 Data

6.1 Feature space

- Measurement space \mathcal{X}
- + numerical $\mathcal{X} \subset \mathbb{R}^d$
- + boolean $\mathcal{X} = \mathbb{B}$
- + categorial $\mathcal{X} = \{1, \dots, k\}$

Features are derived quantities or indirect observations which often significantly compress the information content of measurements.

Remark The selection of a specific feature space predetermines the metric to compare data; this choice is the first significant design decision in a machine learning system.

Taxonomy of Data

6.2 Example of Data

- monadic data
- dyadic data
- pairwise data
- polyadic data

7 Mathematical Spaces

- Topological spaces
- Metric space
- Euclidean vector spaces
- Probability Spaces

Regression

8 Linear Regression

- Statistical model

$$Y = X^T \beta. \quad Y \in \mathbb{R}. \quad X, \beta \in \mathbb{R}^{d+1}$$

- Residual Sum of Squares (RSS)

$$RSS(\beta) = (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta)$$
$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

9 Gauss Markov Theorem

10 Bias/Variance Dilemma

- Tradeoff, split Error
- Identify error components

11 Bayesian Maximum A Posteriori (MAP) estimates

11.1 Ridge Regression

- Cost function
- Bayesian view
- Solution

Tikhonov regularization

11.2 LASSO

- Cost function
- Bayesian view
- Solution

11.3 Ridge vs. LASSO Estimation

12 Remarks on Shrinkage Methods

- Generalized Ridge Regression

Idea behind shrinkage When white noise is added to the data then all Fourier coefficients are increased by a constant on average. ▮ Shrink all coefficients by the estimated noise amount to derive a robust predictor.

13 Learning Objectives

- To motivate, understand, and design Gaussian processes.
- To be able to analytically derive procedures for making predictions with Gaussian processes.
- To analytically compute conditionals, marginals, and posteriors of Gaussians.
- To formulate and understand kernels.
- To be able to use kernel engineering to design new kernels.
- To be able to make a formal connection between Gaussian

processes and Bayesian linear regression.

14 Gaussian Processes

14.1 Bayesian linear regression

multiple linear regression model

$$Y = X^T \beta + \epsilon \quad \text{Gaussian Noise } \epsilon \sim \mathcal{N}(\epsilon|0, \sigma^2)$$

$$p(Y|X, \beta, \sigma) = \mathcal{N}(Y|X^T \beta, \sigma^2) \propto e^{-\frac{1}{2\sigma^2}(Y - X^T \beta)^2}$$

Bayesian linear regression extends multiple linear regression by defining a prior over the regression coefficients, for example (ridge regression)

- Model inversion

14.2 Moments of Bayesian linear regression

Setting

Expected Value

Covariance

15 Gaussian processes

Moments of joint Gaussian:

$$Y \sim \mathcal{N}(Y|0, k_{i,j} + \sigma^2 \text{ if } i = j)$$

with $k_{i,j}$ kernel function

Gaussian Processes as “kernelized linear regression”

- Kernel functions specify the similarity between any two data points.

15.1 Recall

Kernel properties:

- Symmetry

- Positive semi-definit

15.2 Gram matrix

Must be positive semi-definit

$$K = \begin{bmatrix} k(x_1, x_1) & \cdots & k(x_1, x_n) \\ \vdots & & \vdots \\ k(x_n, x_1) & \cdots & k(x_n, x_n) \end{bmatrix}$$

15.3 Examples of kernel functions

Linear kernel: $k(x, x_0) = x^T x_0$

Polynomial kernel: $k(x, x_0) = (x^T x_0 + 1)^p$, for $p \in \mathbb{N}$

Gaussian (RBF) kernel: $k(x, x_0) = \exp -kx - x_0 k_{22} / h^2$

Sigmoid (tanh) kernel: $k(x, x_0) = \tanh kx^T x_0 - b$

Different kernels have different **invariance properties!**

For example, invariance to **rotation** or **translation**.

15.4 Kernel engineering by composition

Addition: Multiplication: Scaling: Composition:

15.5 Prediction by Gaussian processes

Predictive density $p(y_{n+1}|x_{n+1}, X, y)$

Reminder: Conditional Gaussian Distributions

15.6 Prediction by Gaussian processes

15.7 Kernel validation

Goal: Validate hyperparameters of kernels by random splits D

16 Controller Optimization for Robust Control

Machine Learning in Control Systems

Machine learning techniques are becoming more and more important for enabling computers to control complex and stochastic systems and predict the outcomes of such systems.

16.1 Gaussian processes for Control

A Fundamental problem when designing controllers for dynamic systems is the estimation of the controller parameters. Besides pure statistical performance, robustness arises as an important design issue.

The classical approach selects a model of the system to design an initial controller; parameters are then tuned manually to achieve best performance.

An alternative approach uses methods from machine learning to optimize statistical performance, e.g., Bayesian optimization.

Safety-critical system failures may happen because these methods evaluate different controller parameters.

16.2 Safe optimization

Overcome safety-critical system failures by using a specialized optimization algorithm for automatic controller parameter tuning. This algorithm models the underlying performance measure as a GP and only explores new controller parameters whose performance lies above a safe performance threshold with high probability.

17 Model averaging is common practice

- Previous: Gaussian process motivated by Bayesian linear regression.
- Seldom: take MAP estimator in Bayesian setting.
- Bayesian approach: average models with different parameters (weighted according to prior).
- Cross validation: Take average over models trained on different folds.
- Winners of most Machine Learning competitions (e.g. on Kaggle): ensembles (weighted averages of models).

18 Combining Regressors - Bias

TODO: formula

19 Combining Regressors - Variance

TODO: formula

20 Ensemble Learning

The idea of classifier ensembles Boosting is an approach to machine learning based on the idea of creating a highly accurate prediction rule by combining many relatively weak and inaccurate rules.

- Computational advantage
- Statistical advantage

21 Induction Principles for Classifier Selection

- I) Empirical Risk Minimization (ERM) Principle
- II) Bayesian inference by model averaging

22 Motivation for Ensemble Methods

- Train several sufficiently diverse predictors
- Bagging
- Arcing
- Boosting

23 Weak Learners Used for Bagging or Boosting

Combining Classifiers

Bagging Classifiers

Classifier selection: First compare, then bag!

Bagging: The Mechanism

Decision Trees

Random Forests

The Idea of Boosting

AdaBoost

Data Reweighting

Boosted Classifier

Comparison of ensemble methods

24 Loss functions for classification

Support Vector Machines

25 SVM and convex optimization

25.1 Learning objectives

- Lagrange multipliers

- Basics of convex optimization and duality

25.2 Instability of the perceptron

25.3 Maximum-margin criterion

25.4 Lagrange Multipliers

25.5 Convex optimization

25.5.1 The dual

25.5.2 Weak duality and strong duality

25.5.3 Slater's condition

25.5.4 Consequences of strong duality

- The primal optimal solution is the minimizer of the Lagrangian:

- Complementary slackness:

25.5.5 Convex optimization via the dual

25.6 SVMs

- Formula separation band

- Projection

- Normalization trick

- Formalization of primal

Why do we compute the dual?

- SVM calculation can become intractable in high dimensions or when working with transformations (more on this later).

- The dual does not depend on the dimensionality of the dataset in the primal!

Complementary slackness

What about the intercept?

What if the dataset is not linearly separable?

- Soft-margin SVMs

- Kernels

25.7 Soft-margin SVMs

- Formulation

- The role of C

dual ...

26 Kernelization

Transformations and kernels

Polynomial transformations

SVMs and kernels

26.1 Kernels

no need for ϕ only $\phi^T(x)\phi(x')$ for some cases of x and x'

27 SVM summary

Support Vector Machine (SVM): Idea

- margin

- kernel

Nonlinear Transformation in Kernel Space

SVM Lagrangian for functional margin formulation

Non-separable case: Soft Margin SVM

Learning the Soft Margin SVM

28 Structural SVMs

From margins to score functions

Extensions to the SVM

Multiclass SVMs (linear discriminant function)

Structural SVMs

..more examples

Neural Networks

- Activation functions

- Forward computation

- Gradient descent

- Chain rule

A **computation DAG** is a directed acyclic graph where vertices are annotated with variables and edges are annotated with differentiable vector fields.

- generalization chain rule

- Backward computation

Transformer

29 Tokenization and embeddings

- Topic classification

Strategy for NLP

We create for each word in the sentence an embedding vector. Then we aggregate these embeddings to create an embedding for the sentence. We then apply a function on the embedding to map the sentence to a distribution of topics.

- Tokenization

The first step is to transform the sentence into a sequence of tokens.

A popular method is the WordPiece tokenization. It produces a decomposition of a training corpus of text into S tokens, where S is predefined in advance.

1. Set T = set of characters occurring in text

2. While size of $T > S$ do:

1. Find the most common pair a, b of tokens in T occurring in the text.

2. Remove a and b from T and put instead a new token ab in T .

30 Self-attention

- Sentiment classification for movies

- Attention

- Relational matrices

These maps can also be used to infer meanings from words

Every word finds its meaning through other words in the sentence!

Attention as information retrieval

An attention map for location of objects

- Query
- Keys
- Values

Abstraction of an attention mechanism

30.1 Multi-headed attention

30.2 Cross-attention

Consider the problem of translating from English to German:

- A different attention mechanism is needed

Abstraction of a cross-attention mechanism

30.3 Masked self-attention

Consider the task of generating text for answering questions.

31 Positional encodings

- Consider the sentences the dog chased the cat and the cat chased the dog.
- They are permutations of each other, so the attention maps will also just be permutations of each other.
- Note that attention mechanisms do not pay attention to the position of a token in the sentence.
- Binary positional encoding
- Positional encodings

Why adding and not concatenating?

32 Broadcasting

- Broadcasting in PyTorch
- Examples of broadcasting

Broadcasting formalized

Implementing the masked-attention mechanism

33 Residual Blocks

- The degradation problem
- Residual blocks

34 Transformer Architecture

BERT: Bidirectional Encoder Representations from Transformers

Computer Vision

35 CNNs and UNets

35.1 Convolutions and deconvolutions

35.1.1 Convolutional filters

- Input feature map (5x5)
- Padding (2)
- Kernel size (3x3)
- Stride (4)

35.1.2 Deconvolutional filters

35.1.3 Convolutional and deconvolutional layers

- A (de)convolutional layer consists of a number of equally-sized (de)convolutional filters, a padding, and a stride.
- The output is a tensor where the number of channels matches the number of filters in the layer.

35.2 U-net

- A typical architecture for image segmentation.

- One can also treat a U-net as some sort of image autoencoder
- Image from original paper: <https://arxiv.org/abs/1505.04597>

36 Stable diffusion

36.1 Diffusion models with images

Denoising diffusion probabilistic models (DDPM) (Ho et al, 2020)

36.1.1 Forward process

36.1.2 Reverse process

36.2 Conditional diffusion processes

37 Neural networks and variational autoencoder

37.1 Clip

Autoencoders create their own representations

This prevents the decoding as text of image representations produced by an encoder.

How can we harmonize such representations?

38 Stable diffusion

Goal

Given a text prompt, produce an image associated to that prompt.

Idea

Produce text and image embeddings that are semantically relatable.

When given a text, we produce an embedding and then use a generative model to produce an image, guided with the text.

38.1 Semantically relatable text embeddings with CLIP

1. Select an image and a text encoder.
2. Encode all images and texts.
3. Transform all embeddings to the same Euclidean space.
4. Compute a similarity matrix containing all similarities between each pair of text-image embeddings.
5. Maximize the diagonal elements and minimize the non-diagonal elements.

38.2 Cross-attention

We need a cross-attention mechanism for images and texts.

Multi-headed cross-attention mechanism for UNets

39 ViTs and MAEs

39.1 Vision transformers

39.2 Masked auto encoders

39.2.1 Image encoding vs text encoding

Images are signals with very heavy spatial redundancy. Words are signals with low redundancy.

Pixels have little semantic information. Words have high semantic information.

Principle for an image autoencoder

- Remove a large fraction of the image, and create the encoding from there. Then learn to reconstruct the image from only those encodings.

Diffusion

Lecture 8, 03.11.25

Diffusion (to produce images)

- goal find p^*

- Start with gaussian distribution
- generative process, from noise to image
- corruptive process, from image sample to noise

Training

corruptive process

$x(i)$ image

$e_t(i) \rightarrow x_t(i)$

up to the middle / noise

reverse process

from noise to image

Diffusion process for MNIST

39.3 sgrok

Diffusion Models for Image Generation

Diffusion models are a class of generative models aimed at producing images by learning to approximate the true data distribution p^* . The core idea involves two complementary processes: a **corruptive process** that gradually adds noise to an image, transforming it into pure Gaussian noise, and a **reverse process** that denoises the noise back to a realistic image.

Corruptive Process (Forward Diffusion) Starting from a clean image sample $x_0(i)$ (e.g., a digit from the MNIST dataset), the corruptive process applies a sequence of noise additions over T timesteps:

$$x_t(i) = \sqrt{\alpha_t}x_{t-1}(i) + \sqrt{1 - \alpha_t}\epsilon_t(i), \quad \epsilon_t(i) \sim \mathcal{N}(0, I),$$

where α_t controls the noise level, and $\epsilon_t(i)$ is Gaussian noise. This continues until $x_T(i)$ approximates a standard Gaussian distribution (pure noise) at timestep T , effectively destroying the image structure.

Training During training, the model learns to reverse this corruption. A neural network (e.g., a U-Net) is trained to predict the noise $\epsilon_t(i)$ added at each timestep t , given the noisy image $x_t(i)$. The objective is to minimize the difference between predicted and actual noise, often using a mean-squared error loss. This allows the model to capture the data distribution by simulating the forward process on training samples (like MNIST digits) and optimizing the reverse denoising steps.

Generative Process (Reverse Diffusion) For generation, start with a sample from a Gaussian distribution (pure noise) $x_T \sim \mathcal{N}(0, I)$. The trained model iteratively denoises it via the reverse process:

$$x_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \hat{\epsilon}_t(x_t) \right) + \sigma_t z, \quad z \sim \mathcal{N}(0, I)$$

where $\hat{\epsilon}_t$ is the model's noise prediction, and $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$. Over T steps, this transforms noise into a high-fidelity image, such as a synthetic MNIST digit.

This framework, as in Denoising Diffusion Probabilistic Models (DDPM), excels in tasks like image synthesis due to its stable training and high-quality outputs. For implementation, engineers can use libraries like Diffusers (Hugging Face) with pre-trained models for MNIST or extend to advanced variants like Stable Diffusion for conditional generation.

39.4 Encoder

Autoencoder

39.4.1 Clip

- encode image and text
- match them somehow
- produce like diagonal matrix

39.5 Pipeline

Goal: text to image

Idea:

- produce text and image embeddings that are semantically relatable
- text -> produce embedding -> generative model -> image

39.5.1 details

- cross attention

39.6 Unet Architecture

prompt -> encoder -> embedding-text

image -> encoder -> et -> corrupted image xt

xt -> unet (only image input), now?

add text with cross-attention inbetween UNet steps

39.6.1 Multi-headed X Attention Mechanism

embedding-text (Batch,MaxTokens,DimEmbeddings)

embedding-image (Batch,Chanel, Height, Width)

wK,wV to Et, gives $K(B,F,M,D_k), V(B,F,M,D_v)$

wQ to Ei, gives Q

- calculate similarities

use K,Q to create $P(B,F,M,H,W)$

- P = similarity of token m in text and pixel (h,w) of image

create S, softmax of P (along dimension M)

Create $A(B,F,D_v,H,W)$ from S and V

use wO(Do,H,W) to create output $A_o(B,D_o,H,W)$

Graph Neural Networks

40 Graph

Is a molecule toxic or not? structure matters -> graph

Graph $G(V, E)$ consists of

- non-empty vertex set V
- edges $E \subseteq V \times V$

Annotated graph $G = (V, E, h^V, h^E, h_0)$

- $h^V : V \rightarrow \mathbb{R}^{d_v}$ assigns vector of features to each vertex
- $h^E : E \rightarrow \mathbb{R}^{d_E}$ assigns vector of features to each edge
- $h_0 \in \mathbb{R}^{d_0}$ vector of global features

40.1 Neural Networks for graphs

Graphs are not like vectors, who can be represented as arrays of fixed length.

Graphs are not like images, which can be represented as tensors of fixed dimensions.

Graphs are not like texts, which can be represented as sequences of tokens of arbitrary length.

- Standard neural network architectures cannot be used to represent randomized graph functions!

- How can we come up with neural networks that are “universal approximators” of randomized graph functions?

40.1.1 Analyzing how neural networks work

We iteratively process representations by

- transforming them,
- aggregating them with sum,
- and then applying an activation function.

In CNNs, we iteratively process image representations using convolutional filters. Such filters process a patch of the image by:

- transforming the patch
- aggregating their values with an operator (usually sum)
- and then applying an activation function

40.1.2 Generalization to graphs

We then need a “graph filter” that processes “patches” of the graph to produce new representations. Such filter would:

- transform the “patch”,
- aggregate its results,
- and then apply an activation function.

For each vertex, we define a patch as the vertex and its set of neighbors.

- So suppose that the features of each vertex are h_u , for $u \in V$.
- We produce new features as follows:

$$h_u' = \phi \left(\frac{1}{\sqrt{\deg(u) \deg(v)}} \sum_{v \in N(u) \cup u} h_v W \right)$$

- W is a matrix.
- ϕ is an element-wise activation function.
- $\deg(u)$ is the degree of u ; that is, the number of edges adjacent to u .

We divide by the degrees of u and v to maintain all features on a same scale.

We can then iterate this process multiple times, giving rise to layers.

COMPACT MATRIX NOTATION

$$H^{(l+1)} = \phi(X)$$

40.2 Oversmoothing

All representation vectors collapse into a one dimensional subspace.

Solution 1: GIN

Solution 2: Laplace Positional Encodings

Anomaly Detection

41 Data, measurements, and representations

41.1 Four Paradigms in Data Science

41.1.1 Frequentism

1. Pick a parametric model.
 2. Fit the model using maximum likelihood estimation (MLE).
- + Tractable.
 - + Asymptotically unbiased.
 - Variance issues.

41.1.2 Bayesianism

1. Guess a prior.
2. Pick a parametric model.

3. Derive a posterior.

+ Low variance.

- Intractable.

- Bias issues.

41.1.3 Statistical Learning

1. Forget about distributions!

2. Define a loss function.

3. Approximately minimize the expected loss.

+ Tractable.

+ Low bias and variance with proper model.

- Model selection problem.

41.1.4 Non-Parametric Statistics

1. Forget about distributions!

2. Generate multiple samples from your data.

3. Produce estimates from these samples.

+ Low bias and variance.

- Highly intractable.

41.2 Statistical learning

- restrict the space of possible choices of f to a usually parameterizable set H

- collect a sample Z

- use a differentiable loss function

- approximate our desired f with:

41.3 Objects, data, and representations

- represent objects of interest and characterize them according to their typical patterns for detection, classification, abstraction (compression), etc..

- We represent objects using measurements in a sample space.

Features are derived quantities or indirect observations which often significantly compress the information content of measurements.

- Hypothesis class

- Loss function

41.4 Recap

41.4.1 Calculus

- Derivation

- Useful properties

41.4.2 Probability theory

Probability space

Random variables

Expectations

41.4.3 information theory

Bit-calculation, how much do you need to encode ... ?

Code lengths, informativeness, uncertainty, and surprise

Entropy

Conditional entropy

Mutual information

Kullback-Leibler divergence

Cross-entropy

42 Anomaly Detection

Problem formalization

Given $X = \{x_1, \dots, x_n\} \subseteq \mathbb{R}^D$ from *normal* class $N \subseteq \mathbb{R}^D$

Compute $\phi : \mathbb{R}^D \rightarrow \{0, 1\}$ such that $\phi(x) = 1$ iff $x \notin N$

Strategy

- An anomaly is an unlikely event.
- We fit a model of a parametric family of distributions
- We define an anomaly score

The hypothesis class

It has been observed that linear projections of high-dimensional distributions onto low-dimensional spaces resemble Gaussian distributions.

We propose then to project the training data and then fit a GMM to it.

42.1 Dimensionality reduction with PCA

Given $X = \{x_1, \dots, x_n\} \subseteq \mathbb{R}^D$

Find linear projection $\pi : \mathbb{R}^D \rightarrow \mathbb{R}^d$, with $d \ll D$, such that $\pi(X)$ has sufficiently large variance.

Simple case $d = 1$

- Solve it with Lagrange multipliers
- The maximizer of this quadratic function is the largest eigenvalue λ_1^* of S !
- The optimal projection is the one onto a unit eigenvector u_1^* of λ_1^* .

General case $d > 1$

Set $X_0 = X$ and compute forward to d :

- Compute u_{i+1}^* from X_i as before.
- $X_i = \{x - \text{proj}_{u_i^*} x : x \in X_{i-1}\}$

Finally $\pi(x) = (x^T u_1^*, \dots)$

42.1.1 Strategy

- Project the training data X from \mathbb{R}^D to \mathbb{R}^d
- Fit a GMM p_θ to the projected data
- The anomaly score of a point x is then $-\log p_\theta(\pi(x))$

42.2 Gaussian mixture models

A GMM with K components consists of all distributions whose pdf p is of the form:

parameter π_k (weighted norm), μ_k mean Σ_k positive definite variance

42.2.1 Fitting GMMs with maximum likelihood

- choose the distribution p in the model that maximizes the log likelihood
- Decomposition of the log likelihood
- Properties of E and M
- An approximate optimization algorithm

42.3 Anomaly detection with PCA and EM

- compute projector using PCA
- fit GMM? using EM-algorithm
- anomaly score for new point is: $-\log p_\theta(\pi(x))$

42.3.1 Validation metric

F1 score

- only high when both the precision and recall are high

42.3.2 Experimental validation

1. We take several classification datasets. For each dataset, we do the following.

2. For each class, we define the class as normal and all the other

classes as anomalies.

3.The normal class is divided into a training set and a testing set at random.

4.The training set is given as input to our GMM+EM procedure.

5.The trained anomaly detector is then asked to classify all examples in the testing set and the other classes.

6.We measure the F1 score for each class.

7.We average the F1 scores.

Reinforcement Learning

43 (PO-)MDP by construction

A (partially observable) Markov decision process is a tuple

$$(\mathcal{S}, \mathcal{A}, p, r, p_0, \gamma)$$

- \mathcal{S} set of states (measurable space if continuous)

- \mathcal{A} action set (discrete or continuous)

- $p(\cdot | s, a)$ transition kernel over \mathcal{S}

- $r(s, a, s') \in \mathbb{R}$ reward

- p_0 initial-state distribution

- $\gamma \in (0, 1]$ discount?

Observations

In a POMDP, agent perceives $o \sim \mathcal{O}(\cdot | s)$

In a MDP, $o = s$ (full observability)

43.1 Policies

A policy is a function π that maps each state $s \in \mathcal{S}$ to a distribution over the action set $\pi(\cdot | s)$.

We distinguish between two types of policies:

- Deterministic: $\mu : \mathcal{S} \rightarrow \mathcal{A}$, when $\pi(\cdot | s)$ has no randomness.

- Stochastic: $\pi(\cdot | s)$, otherwise.

43.2 Exploration-exploitation

Goal: maximize expected return while reducing uncertainty.

- ϵ -greedy: with probability ϵ pick random action, else $\arg \max_a Q$

- Softmax/Boltzmann: $\pi(a | s) \propto e^{Q(s,a)/\tau}$ with temperature $\tau > 0$

- Anneal ϵ or τ over time.

43.3 Trajectories

Given a policy π , a trajectory is:

$$\tau = (s_0, a_0, r_0, s_1, a_1, r_1, \dots)$$

43.4 Markov property

By construction:

$$\mathbb{P}(s_{t+1} | s_{0:t}, a_{0:t}) = \mathbb{P}(s_{t+1} | s_t, a_t)$$

Consequences:

- Sufficient statistic is current (s_t, a_t)

- Enables dynamic programming / Bellman recursion.

43.5 Value functions

For policy π :

$$V^\pi(s) = \mathbb{E}_\pi[G | s]$$

$$Q^\pi(s, a) = \mathbb{E}_\pi[G | s, a]$$

Relation for discrete \mathcal{A} :

$$V^\pi(s) = \sum_a \pi(a | s) Q^\pi(s, a)$$

43.6 Bellman equations

Policy evaluation:

Optimality (MDP):

43.7 Q-learning (tabular, off-policy)

Update at (s_t, a_t, r_t, s_{t+1}) :

$$Q \leftarrow Q + \alpha[r_t + \gamma \max_{a'} Q(s_{t+1}, a') - Q(s_t, a_t)]$$

43.8 DQN (function approximation)

- Approximate $Q_\theta(s, a)$ with an MLP.

- Target network θ^- ; replay buffer for decorrelation.

- Squared TD error with max over a' using θ^-

44 Active learning

44.1 Information-based transductive learning

44.2 Safe Bayesian optimization

We want to find the maximum of an unknown stochastic process f^\star

We can iteratively choose points $x_1, \dots, x_{n-1} \in X$ and observe realizations of $y_i = f^\star(x_i)$

We must not pick points outside the safe area $S^\star = \{x \in X : g^\star(x) \geq 0\}$, where g^\star is another stochastic process.

We also observe the realizations $z_i = g^\star(x_i)$

How do we pick x_n so that we can estimate best where the maximum in expectation of f^\star is?

44.2.1 Bayesian optimization with ITL

Fit a GPs f on $\{x_i, y_i\}$ for $i < n$

This GP induces 2 functions, l_n^f and u_n^f , defined by requirement:

- $\forall x \in X, [l_n^f(x), u_n^f(x)]$ is 95%-confidence interval of $\mathbb{E}[f(x)]$

Similarly, fit a GP g on $\{x_i, z_i\}$ for $i < n$

This GP induces 2 functions, l_n^g and u_n^g in an analogous way and they in turn induce the following pessimistic and optimistic estimates of the safe set

$$\mathcal{S}_n = \{x : l_n^g(x) \geq 0\} \text{ and } \hat{\mathcal{S}}_n = \{x : u_n^g(x) \geq 0\}$$

We can then consider the following set of estimates

$$\mathcal{A}_n = \{x \in \hat{\mathcal{S}}_n : u_n^f(x) \geq \max_{x' \in \mathcal{S}_n} l_n^f(x')\}$$

We can then do ITL using as sample space \mathcal{S}_n and target space \mathcal{A}_n

44.3 Batch active learning

We assume given the following:

- an input domain \mathcal{X} and a distribution P over \mathcal{X}
- an oracle to an unknown function $f : \mathcal{X} \rightarrow \mathcal{Y}$
- a population set $X = \{x_1, \dots, x_m\} \subseteq \mathcal{X}$
- a budget $b \leq m$

Our goal is to find $L \subseteq X$ such that $|L| = b$ and query the oracle for $\{f(x) : x \in L\}$.

AUXILIARY DEFINITIONS

44.3.1 Problem formalization