# Advanced Machine Learning

**Silvan Stadelmann** - 29. Oktober 2025 - v0.0.1

github.com/silvasta/summary-aml


created with grok

## Contents

# Representation

## 1 Learning objectives

Estimation of Dependences Based on Empirical Data

What is the learning problem?

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

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

## 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 $\mathcal{X} = \{1, ..., k\}$

**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

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

- Residual Sum of Squares (RSS)

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

## 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 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).

-

# 14 Combining Regressors - Bias
TODO: formula

-

# 15 Combining Regressors - Variance
TODO: formula

# 16 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

# 17 Induction Principles for Classifier Selection
I) Empirical Risk Minimization (ERM) Principle

II) Bayesian inference by model averaging

# 18 Motivation for Ensemble Methods
- Train several sufficiently diverse predictors

- Bagging

- Arcing

- Boosting

# 19 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

# 20 Loss functions for classification
# 21 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.

# 22 Gaussian Processes

## 22.1 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

## 22.2 Moments of Bayesian linear regression

Setting

Expected Value

Covariance

# 23 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.

## 23.1 Recall

Kernel properties:

- Symmetry

- Positive semi-definit

## 23.2 Gram matrix

Must be positive semi-definit

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

## 23.3 Examples of kernel functions

Linear kernel: k(x, x0 ) = xT x0

Polynomial kernel: k(x, x0 ) = (xT x0 + 1)p , for p ⬚ N

Gaussian (RBF) kernel: k(x, x0 ) = exp −kx − x0 k22 /h2

Sigmoid (tanh) kernel: k(x, x0 ) = tanh κxT x0 − b

Different kernels have different **invariance properties**!

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

## 23.4 Kernel engineering by composition

Addition: Multiplication: Scaling: Composition:

## 23.5 Prediction by Gaussian processes

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

Reminder: Conditional Gaussian Distributions

## 23.6 Prediction by Gaussian processes

## 23.7 Kernel validation

Goal: Validate hyperparameters of kernels by random splits D

# 24 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.

### 24.1 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.

### 24.2 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.

# Support Vector Machines
# Neural Networks
# Transformer
# Exercises

## 25 Problem 1 - Regression

- Linear Regression

- Ridge Regression

- Noisy Regression

## E1.2.c - An Engineer's rule of thumb is to choose K as $\min \sqrt{n}, 10$

- Overfitting

- Cross Validation

- Generative vs. Discriminative Modeling