
Statistical Machine Learning

Summary

Statistical Machine Learning

May 15, 2020



TECHNISCHE
UNIVERSITÄT
DARMSTADT

Contents

| | |
|---|-----------|
| 1. Introduction | 11 |
| 1.1. Examples | 11 |
| 1.2. Classification vs. Regression | 11 |
| 1.3. Paradigm | 11 |
| 1.4. Key Challenges | 11 |
| 1.4.1. Generalization | 11 |
| 2. Fundamentals: Linear Algebra | 12 |
| 2.1. Vectors | 12 |
| 2.2. Matrices | 13 |
| 2.3. Operations and Linear Transformations | 15 |
| 2.4. Eigenvalues and -vectors | 15 |
| 2.5. Wrap-Up | 16 |
| 3. Fundamentals: Statistics | 17 |
| 3.1. Random Variables | 17 |
| 3.2. Distributions | 17 |
| 3.2.1. Uniform Distribution | 17 |
| 3.2.2. Discrete Distributions | 17 |
| 3.2.3. Continuous Distributions | 19 |
| 3.2.4. Multivariate Gaussian Distribution | 20 |
| 3.2.5. Partitioned Gaussian Distributions | 21 |
| 3.3. Central Limit Theorem | 21 |
| 3.4. Probability Rules | 21 |
| 3.5. Expectation, Variance and Moments | 22 |
| 3.5.1. Expectation | 22 |
| 3.5.2. Variance and Covariance | 22 |
| 3.5.3. Moments | 23 |
| 3.6. Exponential Family | 23 |
| 3.6.1. Example: Bernoulli Distribution | 23 |
| 3.6.2. Example: Gaussian Distribution | 23 |
| 3.7. Information Theory and Entropy | 24 |
| 3.7.1. Information and Entropy | 24 |
| 3.7.2. Kullback-Leibler Divergence | 24 |
| 3.8. Wrap-Up | 24 |
| 4. Fundamentals: Optimization | 25 |
| 4.1. Convexity | 25 |
| 4.2. Cost Functions | 25 |
| 4.2.1. Common Cost Functions | 25 |
| 4.3. Constrained/Unconstrained Optimization | 25 |
| 4.4. Lagrange Multipliers | 26 |
| 4.4.1. Dual Formulation | 26 |
| 4.4.2. Example | 27 |
| 4.5. Numerical Optimization | 27 |
| 4.5.1. Learning Rate | 28 |
| 4.5.2. Test Functions | 28 |
| 4.5.3. Axial Iteration | 28 |
| 4.5.4. Steepest Descent | 30 |
| 4.5.5. Newtons Method | 30 |
| 4.5.6. Quasi-Newton Method (BFGS) | 35 |
| 4.5.7. Conjugate Gradient (CG) | 35 |
| 4.5.8. Conjugate Gradients vs. BFGS | 38 |

| | |
|--|-----------|
| 4.6. Wrap-Up | 38 |
| 5. Bayesian Decision Theory | 41 |
| 5.1. Character Recognition | 41 |
| 5.1.1. Class Conditional Probabilities | 41 |
| 5.1.2. Class Priors | 41 |
| 5.2. Bayesian Decision Theory | 43 |
| 5.3. Bayesian Probabilities | 43 |
| 5.4. Misclassification Rate | 43 |
| 5.5. Decision Rule, Optimal Classifier and Decision Boundary | 43 |
| 5.5.1. Multiple Classes | 43 |
| 5.5.2. High Dimensional Features | 43 |
| 5.6. Dummy Classes | 43 |
| 5.7. Risk Minimization | 44 |
| 5.7.1. Decision Rule | 44 |
| 5.8. Wrap-Up | 44 |
| 6. Probability Density Estimation | 45 |
| 6.1. Parametric Models | 45 |
| 6.1.1. Maximum Likelihood | 45 |
| 6.1.2. Degenerate Case | 46 |
| 6.1.3. Bayesian Estimation | 47 |
| 6.2. Non-Parametric Models | 48 |
| 6.2.1. Histograms | 48 |
| 6.2.2. Kernel Density Estimation (KDE) | 48 |
| 6.2.3. K-Nearest Neighbors (KNN) | 52 |
| 6.3. Mixture Models | 54 |
| 6.3.1. Mixture of Gaussians | 54 |
| 6.3.2. Estimation using Clustering | 56 |
| 6.3.3. Mixture Components | 60 |
| 6.4. Wrap-Up | 60 |
| 7. Clustering | 61 |
| 7.1. Mean Shift Clustering | 61 |
| 7.2. Wrap-Up | 65 |
| 8. Evaluation | 66 |
| 8.1. Test Error vs. Training Error | 66 |
| 8.2. Bias and Variance | 66 |
| 8.2.1. MVUE and BLUE | 66 |
| 8.2.2. Bias-Variance Tradeoff | 66 |
| 8.2.3. Example: MLE of a Gaussian | 66 |
| 8.2.4. Example: Regression | 68 |
| 8.3. Model Selection and Occam's Razor | 68 |
| 8.3.1. Cross Validation | 68 |
| 8.3.2. K -Fold Cross Validation | 69 |
| 8.3.3. Machine Learning Cycle | 69 |
| 8.4. Wrap-Up | 69 |
| 9. Regression | 70 |
| 9.1. Linear Regression | 70 |
| 9.1.1. Least Squares Regression | 70 |
| 9.2. Generalized Linear Regression | 72 |
| 9.3. Maximum Likelihood Approach | 72 |
| 9.3.1. Probabilistic Regression | 72 |
| 9.3.2. Maximum Likelihood Regression | 74 |
| 9.3.3. Loss Functions | 74 |
| 9.4. Bayesian Linear Regression | 78 |
| 9.4.1. Maximum A-Posteriori (MAP) | 78 |
| 9.4.2. Full Bayesian Regression | 78 |

| | |
|---|------------|
| 9.5. Kernel Regression | 79 |
| 9.5.1. Dual Representation of Regression | 85 |
| 9.5.2. Useful Kernels | 85 |
| 9.6. Gaussian Processes Regression | 86 |
| 9.6.1. Regression | 86 |
| 9.6.2. Function Value Prediction | 86 |
| 9.6.3. Conclusion | 87 |
| 9.7. Wrap-Up | 87 |
| 10. Classification | 89 |
| 10.1. Generative vs. Discriminative | 89 |
| 10.2. Discriminant Functions | 89 |
| 10.2.1. Multiple Classes | 89 |
| 10.2.2. Linear Discriminant Functions | 91 |
| 10.3. Fisher Discriminant Analysis | 91 |
| 10.3.1. Least Squares Classification | 91 |
| 10.3.2. Fishers' Linear Discriminant | 92 |
| 10.4. Perceptron Algorithm | 95 |
| 10.4.1. Intuition | 95 |
| 10.4.2. Linear Separability | 95 |
| 10.5. Probabilistic Discriminative Models | 95 |
| 10.5.1. Logistic Regression | 97 |
| 10.6. Wrap-Up | 97 |
| 11. Linear Dimensionality Reduction | 98 |
| 11.1. Introduction | 98 |
| 11.2. Principal Component Analysis | 99 |
| 11.2.1. Derivation | 99 |
| 11.2.2. Conclusion | 100 |
| 11.3. Choosing the target Dimension | 100 |
| 11.4. Applications | 103 |
| 11.5. Wrap-Up | 103 |
| 12. Statistical Learning Theory | 104 |
| 12.1. Supervised Learning | 104 |
| 12.2. Assessment of Optimality: Risk | 104 |
| 12.2.1. Empirical vs. True Risk | 105 |
| 12.2.2. Convergence Properties | 105 |
| 12.3. Risk Bound | 105 |
| 12.3.1. VC-Dimension | 106 |
| 12.3.2. Example | 106 |
| 12.4. Structural Risk Minimization | 106 |
| 12.5. Wrap-Up | 106 |
| 13. Neural Networks | 107 |
| 13.1. Abstraction of a Neuron | 107 |
| 13.2. Single-Layer Neural Networks | 108 |
| 13.2.1. Logistic Regression | 108 |
| 13.2.2. Multi-Class Network | 108 |
| 13.2.3. Least-Squares Loss Function | 108 |
| 13.2.4. Learning with Gradient Descent | 108 |
| 13.3. Multi-Layer Neural Networks | 109 |
| 13.3.1. One hidden Layer? | 109 |
| 13.3.2. Model Type and Model Class | 109 |
| 13.4. Output Neurons, Activation and Loss Functions | 110 |
| 13.4.1. Output Neurons | 110 |
| 13.4.2. Loss Functions | 110 |
| 13.4.3. Activation Functions | 110 |
| 13.5. Forward- and Backpropagation | 112 |
| 13.5.1. Backpropagation | 112 |
| 13.5.2. Formulas | 113 |

| | |
|---|------------|
| 13.5.3. Approximating the Gradient | 113 |
| 13.6. Gradient Descent | 114 |
| 13.6.1. When to update W ? | 114 |
| 13.6.2. Adaptive Learning Rate | 114 |
| 13.6.3. Small Neural Networks | 115 |
| 13.6.4. Initialization | 116 |
| 13.7. Overfitting | 117 |
| 13.7.1. Batch Normalization | 117 |
| 13.8. Theoretical Results | 117 |
| 13.9. Other Network Architectures | 118 |
| 13.9.1. Convolutional Neural Network (CNN) | 118 |
| 13.9.2. Recurrent Neural Network (RNN) | 118 |
| 13.9.3. Long Short-Term Memory Network (LSTM) | 118 |
| 13.10. Applications | 119 |
| 13.10.1. Computer Vision | 119 |
| 13.10.2. Autonomous Systems | 119 |
| 13.11. Radial Basis Function Networks | 119 |
| 13.12. Wrap-Up | 120 |
| 14. Support Vector Machines | 121 |
| 14.1. Linear SVMs | 121 |
| 14.1.1. Optimization Formulation | 121 |
| 14.1.2. Sparsity | 124 |
| 14.2. Nonlinear SVMs | 124 |
| 14.2.1. Optimization Formulation | 124 |
| 14.2.2. Kernel Trick | 125 |
| 14.3. Non-Separable Data | 126 |
| 14.3.1. Slack Variables | 126 |
| 14.3.2. Lack of Sparseness | 127 |
| 14.4. Applications | 127 |
| 14.4.1. Text Classification | 127 |
| 14.4.2. Handwritten Digit Classification | 127 |
| 14.4.3. Support Vector Regression | 127 |
| 14.5. Wrap-Up | 127 |
| A. Self-Test Questions | 128 |
| A.1. Demo | 128 |
| A.2. Organization | 129 |
| A.3. Linear Algebra Refresher | 129 |
| A.4. Statistics Refresher | 130 |
| A.5. Optimization Refresher | 131 |
| A.6. Bayesian Decision Theory | 132 |
| A.7. Probability Density Estimation | 132 |
| A.8. Clustering and Evaluation | 133 |
| A.9. Regression | 134 |
| A.10. Classification | 134 |
| A.11. Linear Dimensionality Reduction and Statistical Learning Theory | 135 |
| A.12. Neural Networks | 136 |
| A.13. Support Vector Machines | 138 |
| A.14. Kernel Regression and Gaussian Processes | 139 |
| B. Code | 140 |
| B.1. Utility | 140 |
| B.1.1. <code>genData.py</code> | 140 |
| B.2. Optimization | 142 |
| B.2.1. <code>optimization.py</code> | 142 |
| B.3. Probability Density Estimation | 145 |
| B.3.1. <code>nonParametricModels.py</code> | 145 |
| B.3.2. <code>mixtureModels.py</code> | 146 |
| B.4. Clustering | 147 |
| B.4.1. <code>clustering.py</code> | 147 |

| | |
|---|-----|
| B.5. Regression | 149 |
| B.5.1. <code>regression.py</code> | 149 |
| B.6. Classification | 151 |
| B.6.1. <code>classification.py</code> | 151 |
| B.7. Linear Dimensionality Reduction | 152 |
| B.7.1. <code>pca.py</code> | 152 |
| B.8. Support Vector Machines | 153 |
| B.8.1. <code>svm.py</code> | 153 |

List of Figures

| | |
|---|-----|
| 2.1. Illustration of Vector Projection | 13 |
| 3.1. Uniform Distribution | 17 |
| 3.2. Binomial Distribution $\text{Bin}(m 10, 0.25)$ | 18 |
| 3.3. Poisson Distribution $p(m 5)$ | 20 |
| 3.4. Standard Gaussian Distribution $\mathcal{N}(x 0, 1)$ | 21 |
| 4.1. Quadratic Function | 29 |
| 4.2. Rosenbrock Function | 29 |
| 4.3. Steepest Descent on Rosenbrock | 31 |
| 4.4. Steepest Descent on Quadratic | 32 |
| 4.5. Newtons Method on Rosenbrock | 33 |
| 4.6. Newtons Method on Quadratic | 34 |
| 4.7. BFGS on Rosenbrock | 36 |
| 4.8. BFGS on Quadratic | 37 |
| 4.9. CG on Rosenbrock | 39 |
| 4.10. CG on Quadratic | 40 |
| 5.1. Class Conditional Probabilities | 42 |
| 5.2. Class Priors | 42 |
| 6.1. Histogram | 49 |
| 6.2. Kernel Density Estimation (Parzen Window) | 50 |
| 6.3. Kernel Density Estimation (Gaussian Kernel) | 51 |
| 6.4. K-Nearest Neighbors | 53 |
| 6.5. Mixture of Gaussians | 55 |
| 6.6. EM for Univariate Gaussian | 58 |
| 7.1. Cluster Gaussians | 62 |
| 7.2. Mean Shift Clustering | 63 |
| 7.3. Mean Shift Clustering (Way) | 64 |
| 8.1. Machine Learning Cycle | 69 |
| 9.1. Regression: True Function | 71 |
| 9.2. Regression: Least Squares, Underfitting | 73 |
| 9.3. Regression: Maximum Likelihood, Underfitting | 75 |
| 9.4. Regression: Maximum Likelihood, Just Right | 76 |
| 9.5. Regression: Maximum Likelihood, Overfitting | 77 |
| 9.6. Regression: Full Bayesian Regression (2 Samples) | 80 |
| 9.7. Regression: Full Bayesian Regression (8 Samples) | 81 |
| 9.8. Regression: Full Bayesian Regression (All Samples) | 82 |
| 9.9. Regression: Kernel Regression (RBF, $\sigma^2 = 0.01$) | 83 |
| 9.10. Regression: Kernel Regression (RBF, $\sigma^2 = 1$) | 84 |
| 10.1. Classification: Example Data (Linear Separable) | 90 |
| 10.2. Linear Separability | 91 |
| 10.3. Classification: Least Squares | 93 |
| 10.4. Classification: Perceptron (7 Iterations) | 96 |
| 11.1. Principal Component Analysis: Iris Dataset | 101 |
| 11.2. Principal Component Analysis: Iris Dataset (Explained Variance) | 102 |

| | |
|--|-----|
| 13.1.Neural Network: Single-Layer | 108 |
| 13.2.Neural Network: Multi-Layer | 109 |
| 13.3.Sigmoid $\sigma(z)$ | 111 |
| 13.4.Sigmoid $\tanh(z)$ | 111 |
| 13.5.Rectified Linear Unit ReLU $\max(0, z)$ | 112 |
| 14.1.Linear Support Vector Machine | 122 |



List of Tables

4.1. Common Cost Functions 26

List of Algorithms

| | | |
|----|--|----|
| 1. | Steepest Descent (Minimization) | 30 |
| 2. | Newtons Method (Minimization) | 30 |
| 3. | Quasi-Newton-Method, BFGS (Minimization) | 35 |
| 4. | Conjugate Gradients (Minimization) | 38 |
| 5. | EM for Univariate Gaussian | 57 |
| 6. | Mean Shift Clustering | 61 |
| 7. | Perceptron Algorithm | 95 |

1. Introduction

Most of the content in this summary, the ideas, the underlying structure and the image ideas are taken from the lecture "Statistical Machine Learning" by Prof. Jan Peters. It is really just a *summary* of the contents of the lecture.

1.1. Examples

1.2. Classification vs. Regression

1.3. Paradigm

1.4. Key Challenges

1.4.1. Generalization

2. Fundamentals: Linear Algebra

2.1. Vectors

A *vector* is an ordered list of numbers that can be interpreted as an *arrow* in multidimensional spaces:

$$\mathbf{v} \in \mathbb{R}^n \rightarrow v = \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix} \quad (2.1)$$

Scalar Multiplication Multiplying a vector by a scalar is defined as multiplying each component by that scalar (let $\mathbf{v} \in \mathbb{R}^n$, $\lambda \in \mathbb{R}$):

$$\lambda \mathbf{v} = \lambda \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix} = \begin{bmatrix} \lambda v_1 \\ \vdots \\ \lambda v_n \end{bmatrix} \quad (2.2)$$

Per component:

$$(\lambda v)_i = \lambda v_i \quad (2.3)$$

Scalar multiplication is a linear operation.

Addition Adding two vectors is defined by adding the component of both vectors (thus, both vectors must have the same size; let $v, w \in \mathbb{R}^n$):

$$\mathbf{v} + \mathbf{w} = \begin{bmatrix} v_1 + w_1 \\ \vdots \\ v_n + w_n \end{bmatrix} \quad (2.4)$$

Per component:

$$(v + w)_i = v_i + w_i \quad (2.5)$$

Vector addition is both associative and commutative.

Vector Transpose The *transposed* version \mathbf{v}^T of a vector $\mathbf{v} \in \mathbb{R}^n$ is a vector that was flipped around its main axis (thus, the new vector is a row/column vector if the initial vector as a column/row vector):

$$\begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix}^T = [v_1 \quad \cdots \quad v_n] \quad (2.6)$$

$$[v_1 \quad \cdots \quad v_n]^T = \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix} \quad (2.7)$$

Transposing a vector twice returns the initial vector

$$\mathbf{v} = (\mathbf{v}^T)^T$$

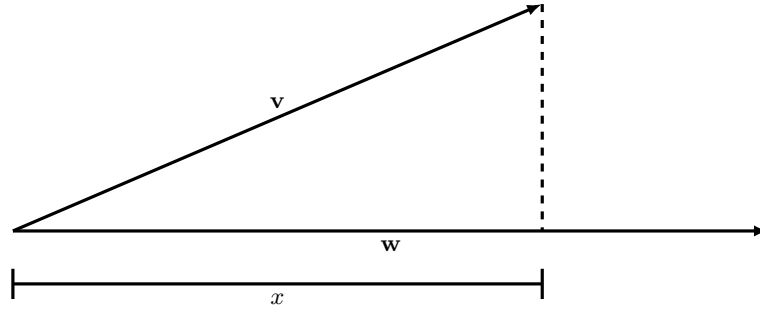


Figure 2.1.: Illustration of Vector Projection

Linear Combination A *linear combination* of multiple vectors $\mathbf{v}_1, \dots, \mathbf{v}_n \in \mathbb{R}^m$ is the addition of the scaled versions of them (scaled by scalars $\lambda_1, \dots, \lambda_n \in \mathbb{R}$):

$$\mathbf{u} = \lambda_1 \mathbf{v}_1 + \dots + \lambda_n \mathbf{v}_n \quad (2.8)$$

If in a group of vectors $\mathbf{v}_1, \dots, \mathbf{v}_n \in \mathbb{R}^m$, no vector can be represented as a linear combination of the others, they are called *linearly independent*.

Inner and Outer Product and Length The *inner product* of two vectors $\mathbf{v}, \mathbf{w} \in \mathbb{R}^n$ is the sum of the product of the components and is denoted by a single dot ($\mathbf{v} \cdot \mathbf{w}$):

$$\mathbf{v} \cdot \mathbf{w} = \mathbf{v}^T \mathbf{w} = (v_1 w_1) + \dots + (v_n w_n) \quad (2.9)$$

Therefore, the inner product gives a scalar value. In Cartesian coordinates, this is also called the *scalar product*.

The length $\|\mathbf{v}\|$ of a vector $\mathbf{v} \in \mathbb{R}^n$ is given as the euclidean norm:

$$\|\mathbf{v}\| = (\mathbf{v} \cdot \mathbf{v})^{\frac{1}{2}} = \sqrt{v_1^2 + \dots + v_n^2}$$

The *outer product* of two vectors $\mathbf{v}, \mathbf{w} \in \mathbb{R}^n$ is defined analogue to the inner product, but with the transpose switched and is denoted by a cross in a circle ($\mathbf{v} \otimes \mathbf{w}$):

$$\mathbf{v} \otimes \mathbf{w} = \mathbf{v} \otimes \mathbf{w}^T \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix} \otimes [w_1 \quad \dots \quad w_n] = \begin{bmatrix} v_1 w_1 & \dots & v_1 w_n \\ \vdots & \ddots & \vdots \\ v_n w_1 & \dots & v_n w_n \end{bmatrix} \quad (2.10)$$

This is equivalent to matrix multiplication with two vectors and thus produces a matrix $\mathbf{v} \otimes \mathbf{w} \in \mathbb{R}^{n \times n}$.

Angles between Vectors The *angle* θ between two vectors $\mathbf{v}, \mathbf{w} \in \mathbb{R}$ is given by:

$$\cos \theta = \frac{\mathbf{v} \cdot \mathbf{w}}{\|\mathbf{v}\| \|\mathbf{w}\|} \quad \Longleftrightarrow \quad \theta = \arccos \frac{\mathbf{v} \cdot \mathbf{w}}{\|\mathbf{v}\| \|\mathbf{w}\|} \quad (2.11)$$

Projections of Vectors A *projection* of a vector $\mathbf{v} \in \mathbb{R}^n$ onto a vector $\mathbf{w} \in \mathbb{R}^n$ results in a scalar value $x \in \mathbb{R}$ which equals the length of the adjacent w.r.t. to the angle between both vectors given that \mathbf{v} is the hypotenuse and the third line is orthogonal to \mathbf{w} , see figure 2.1 for an illustration. Then this length is given as:

$$x = \|\mathbf{v}\| \cos \theta = \frac{\mathbf{v} \cdot \mathbf{w}}{\|\mathbf{w}\|} \quad (2.12)$$

2.2. Matrices

A *matrix* is an ordered group of numbers that are ordered in two dimensions:

$$A \in \mathbb{R}^{n \times m} \quad \rightarrow \quad A = \begin{bmatrix} a_{11} & \dots & a_{1m} \\ \vdots & \ddots & \vdots \\ a_{n1} & \dots & a_{nm} \end{bmatrix} \quad (2.13)$$

Scalar Multiplication Multiplying a matrix by a scalar is defined as multiplying each component by that scalar (let $A \in \mathbb{R}^{n \times m}$, $\lambda \in \mathbb{R}$):

$$\lambda A = \lambda \begin{bmatrix} a_{11} & \cdots & a_{1m} \\ \vdots & \ddots & \vdots \\ a_{n1} & \cdots & a_{nm} \end{bmatrix} = \begin{bmatrix} \lambda a_{11} & \cdots & \lambda a_{1m} \\ \vdots & \ddots & \vdots \\ \lambda a_{n1} & \cdots & \lambda a_{nm} \end{bmatrix} \quad (2.14)$$

Per component:

$$(\lambda A)_{ij} = \lambda a_{ij} \quad (2.15)$$

Scalar multiplication is a linear operation.

Addition Adding two matrices is defined by adding the component of both matrices (thus, both matrices must have the same size; let $A, B \in \mathbb{R}^{n \times m}$):

$$A + B = \begin{bmatrix} a_{11} + b_{11} & \cdots & a_{1m} + b_{1m} \\ \vdots & \ddots & \vdots \\ a_{n1} + b_{n1} & \cdots & a_{nm} + b_{nm} \end{bmatrix} \quad (2.16)$$

Per component:

$$(A + B)_{ij} = a_{ij} + b_{ij} \quad (2.17)$$

Matrix addition is both associative and commutative.

Transpose The *transposed* version A^T of a matrix $A \in \mathbb{R}^{n \times m}$ is a matrix $A^T \in \mathbb{R}^{m \times n}$ that was flipped around its main axis:

$$\begin{bmatrix} a_{11} & \cdots & a_{1m} \\ \vdots & \ddots & \vdots \\ a_{n1} & \cdots & a_{nm} \end{bmatrix}^T = \begin{bmatrix} a_{11} & \cdots & a_{n1} \\ \vdots & \ddots & \vdots \\ a_{1m} & \cdots & a_{nm} \end{bmatrix} \quad (2.18)$$

Transposing a matrix twice returns the initial matrix

$$A = (A^T)^T$$

Matrix Multiplication The *matrix multiplication* of two matrices is only possible if the number of columns of the first matrix equals the number of rows of the second matrix (i.e. $A \in \mathbb{R}^{n \times m}$, $B \in \mathbb{R}^{m \times o}$). The resulting matrix then has the dimensions $AB \in \mathbb{R}^{n \times o}$. Matrix multiplication is defined as follows:

$$AB = \begin{bmatrix} a_{11} & \cdots & a_{1m} \\ \vdots & \ddots & \vdots \\ a_{n1} & \cdots & a_{nm} \end{bmatrix} \cdot \begin{bmatrix} b_{11} & \cdots & b_{1o} \\ \vdots & \ddots & \vdots \\ b_{m1} & \cdots & b_{mo} \end{bmatrix} = \begin{bmatrix} a_{11}b_{11} + \cdots + a_{1m}b_{m1} & \cdots & a_{11}b_{1o} + \cdots + a_{1m}b_{mo} \\ \vdots & \ddots & \vdots \\ a_{n1}b_{11} + \cdots + a_{nm}b_{m1} & \cdots & a_{n1}b_{1o} + \cdots + a_{nm}b_{mo} \end{bmatrix} \quad (2.19)$$

Per component:

$$(AB)_{ij} = \sum_{k=1}^m A_{ik} B_{kj} \quad (2.20)$$

Matrix multiplication is only associative and distributive w.r.t. matrix addition.

Inverse Let $I_n \in \mathbb{R}^{n \times n}$ be the identity matrix with all ones on the main diagonal and the rest zeros. If the dimension is clear, the n can be left out.

Using this definition, the *inverse* of a matrix $A \in \mathbb{R}^{n \times n}$ is defined as matrix $A^{-1} \in \mathbb{R}^{n \times n}$ that holds the following equation:

$$AA^{-1} = A^{-1}A = I_n \quad (2.21)$$

If such a matrix exists, the matrix A is called *regular* or *nonsingular*.

Pseudoinverse If a matrix $A \in \mathbb{R}^{n \times m}$ is not squared (i.e. $n \neq m$), there exists no inverse matrix. Instead, *Pseudoinverse Matrices* can be used to “invert” such a matrix. Left and right pseudoinverse are mutually exclusive, meaning that the one can only exist if the other does not (whilst neither have to exist).

The *left pseudoinverse* does only exist if the matrix has full column rank and is defined as:

$$A^\# = (A^T A)^{-1} A^T \implies A^\# A = I_m \quad (2.22)$$

The *right pseudoinverse* does only exist if the matrix has full row rank and is defined as:

$$A^\# = J^T (J J^T)^{-1} \implies A A^\# = I_n \quad (2.23)$$

Properties

Symmetry A squared matrix $A \in \mathbb{R}^{n \times n}$ is *symmetric* iff $A^T = A$. This implies that:

- The inverse matrix A^{-1} is also symmetric.
- A can be decomposed into $A = Q D Q^T$, where D is a diagonal matrix with all eigenvalues of A and Q is a matrix with all columns as the eigenvectors of A .

Definite Quadratic Form Let $A \in \mathbb{R}^{n \times n}$ be a squared matrix and let $\sigma(A)$ be its spectral. Then the matrix A is

$$\begin{cases} \text{positive definite} & \forall \lambda \in \sigma : \lambda > 0 & \iff & \mathbf{x}^T A \mathbf{x} > 0 \\ \text{negative definite} & \forall \lambda \in \sigma : \lambda < 0 & \iff & \mathbf{x}^T A \mathbf{x} < 0 \\ \text{positive semi-definite} & \forall \lambda \in \sigma : \lambda \geq 0 & \iff & \mathbf{x}^T A \mathbf{x} \geq 0 \\ \text{positive semi-definite} & \forall \lambda \in \sigma : \lambda \leq 0 & \iff & \mathbf{x}^T A \mathbf{x} \leq 0 \\ \text{indefinite} & \text{else} & & \end{cases} \quad (2.24)$$

for all vectors $\mathbf{x} \in \mathbb{R}^n$.

Regularity/Nonsingularity All of the following are equivalent w.r.t. a matrix $A \in \mathbb{R}^{n \times n}$:

- The matrix is regular.
- The matrix is nonsingular (or not singular).
- There exists a matrix $A^{-1} \in \mathbb{R}^{n \times n}$ with $A A^{-1} = A^{-1} A = I_n$.
- The determinant of the matrix is nonzero: $\det A \neq 0$.
- The matrix has full row rank.
- The matrix has full column rank.

2.3. Operations and Linear Transformations

Change of Basis

Linear Transformations

2.4. Eigenvalues and -vectors

Basis

Linear Transformations

2.5. Wrap-Up

- Vectors and matrices
- Operations on vectors and matrices
- Eigenvectors and -values
- Linear transformations

3. Fundamentals: Statistics

For a much more detailed introduction into the basic concepts (e.g. random variables), please lookup the summary of “Math 3: Stochastic and Statistics” (see <https://www.dmken.com/cs>), but notice that it is in German.

3.1. Random Variables

A *random variable* is a number that is determined by chance, draw according to a probability distribution.

3.2. Distributions

A probability distribution describes the probability that a random variable will equal a certain value (or lie in a certain range).

3.2.1. Uniform Distribution

All data/all values are equally likely within a bounded region R with size R .

$$p(x) = \frac{1}{R} \quad (3.1)$$

The distribution is plotted in figure 3.1.

3.2.2. Discrete Distributions

The random variables take *discrete values* (can be infinite, but countably infinity) and their probabilities sum up to 1:

$$\sum_i p(x_i) = 1 \quad (3.2)$$

A discrete distribution is described by a *probability mass function* which is a normalized histogram.

Bernoulli Distribution

A *Bernoulli random variable* only takes on two values, e.g. 0 or 1.

Parameters

μ The probability that the variable equals 1.

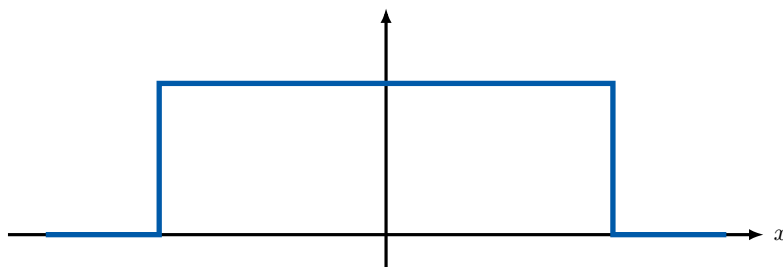


Figure 3.1.: Uniform Distribution

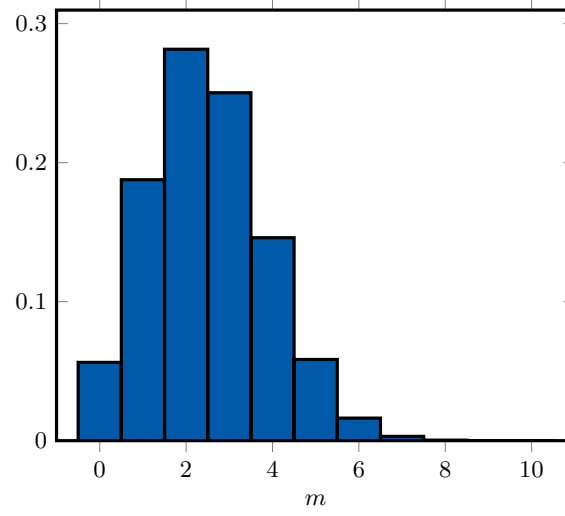


Figure 3.2.: Binomial Distribution $\text{Bin}(m \mid 10, 0.25)$

Properties

$$x \in \{0, 1\} \quad (3.3)$$

$$p(x = 1 \mid \mu) = \mu \quad (3.4)$$

$$\text{Bern}(x \mid \mu) = \mu^x (1 - \mu)^{1-x} \quad (3.5)$$

$$\mathbb{E}(x) = \mu \quad (3.6)$$

$$\text{Var}(x) = \mu(1 - \mu) \quad (3.7)$$

Binomial Distribution

Binomial variables are a sequence of N Bernoulli variables.

Parameters

μ The probability that one variable equals 1.

N The number of trials/samples.

Properties

$$\text{Bin}(m \mid N, \mu) = \binom{N}{m} \mu^m (1 - \mu)^{N-m} \quad (3.8)$$

$$\mathbb{E}(m) = N\mu \quad (3.9)$$

$$\text{Var}(m) = N\mu(1 - \mu) \quad (3.10)$$

See figure 3.2 for a visualization of $\text{Bin}(m \mid 10, 0.25)$.

Multinoulli Distribution

Multinoulli variables (also called *categorical variables*) are a generalization of Bernoulli variables where each variable can have multiple (namely K) outputs. The random variables is a vector with one-hot-encoding.

Parameters

μ The entry μ_i defines the probability that the entry x_i equals 1.
All entries must be $\mu_i \geq 0$ and $\sum_{k=1}^K \mu_k = 1$.

K The number of classes/outcomes.

Properties

$$\mathbf{x} = [0, 0, 1, 0, 0, 0]^T \quad (3.11)$$

$$p(x_i | \boldsymbol{\mu}) = \mu_i \quad (3.12)$$

$$p(\mathbf{x} | \boldsymbol{\mu}) = \prod_{k=1}^K \mu_k^{x_k} \quad (3.13)$$

$$\mathbb{E}(\mathbf{x} | \boldsymbol{\mu}) = \sum_{\mathbf{x}} p(\mathbf{x} | \boldsymbol{\mu}) \mathbf{x} = \boldsymbol{\mu}^T \quad (3.14)$$

Multinomial Distribution

Multinomial variables are a sequence of N Multinoulli variables.

Parameters

$\boldsymbol{\mu}$ The entry μ_i defines the probability that, for one variable, the entry x_i equals 1.
All entries must be $\mu_i \geq 0$ and $\sum_{k=1}^K \mu_k = 1$.

K The number of classes/outcomes.

N The number of trials/samples.

Properties

$$\text{Mult}(m_1, m_2, \dots, m_K | \boldsymbol{\mu}, N) = \binom{N}{m_1, m_2, \dots, m_K} \prod_{k=1}^K \mu_k^{m_k} \quad (3.15)$$

$$\mathbb{E}(m_k) = N\mu_k \quad (3.16)$$

$$\text{Var}(m_k) = N\mu_k(1 - \mu_k) \quad (3.17)$$

$$\text{Cov}(m_j, m_k) = -N\mu_j\mu_k \quad (3.18)$$

Poisson Distribution

A *Poisson distribution* is a binomial distribution where the number of trials goes to infinity $N \rightarrow \infty$ and the success of each trial goes to zero $\mu \rightarrow 0$, s.t. $N\mu = \lambda$ is constant.

Parameters

λ Defines the expectation value and the variance at once.

Properties

$$p(m | \lambda) = \frac{\lambda^m}{m!} e^{-\lambda} \quad (3.19)$$

$$\mathbb{E}(m) = \lambda \quad (3.20)$$

$$\text{Var}(m) = \lambda \quad (3.21)$$

See figure 3.3 for a visualization of $p(m | 5)$.

3.2.3. Continuous Distributions

The random variables take *discrete values* (infinite, can be uncountable) and their probability density function integrates to 1:

$$\int_{-\infty}^{+\infty} p(x) dx = 1 \quad (3.22)$$

A continuous distribution is described by a *probability density function* $p(x)$.

The probability that a random variable x falls into the interval (a, b) is

$$P(a < x < b) = \int_a^b p(x) dx \quad (3.23)$$

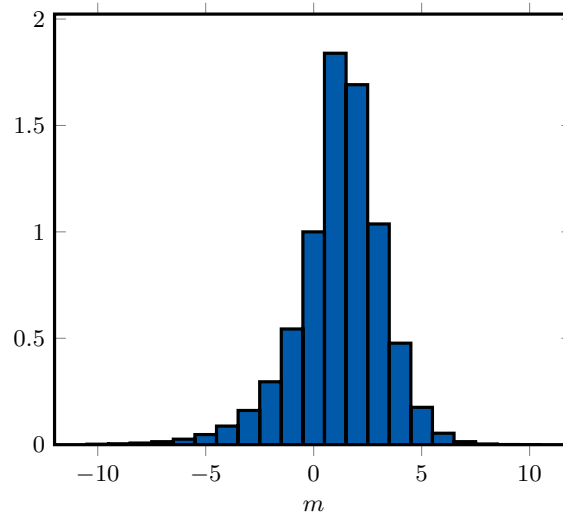


Figure 3.3.: Poisson Distribution $p(m | 5)$

Gaussian Distribution

Parameters

μ The expectation value.

σ^2 The variance.

Properties

$$p(x) = \mathcal{N}(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{1}{2\sigma^2} (x - \mu)^2 \right\} \quad (3.24)$$

$$\mathbb{E}(x) = \mu \quad (3.25)$$

$$\text{Var}(x) = \sigma^2 \quad (3.26)$$

A Gaussian distribution has more really useful properties:

- A Gaussian has soft tails, i.e. they fade away smoothly.
- Gaussians are often good models for data and provide analytical solutions.

See figure 3.4 for a visualization of $\mathcal{N}(x | 0, 1)$ (the standard Gaussian distribution).

3.2.4. Multivariate Gaussian Distribution

Gaussians can be applied to D -dimensional data x_1, x_2, \dots using multivariate Gaussian distributions.

Parameters

μ A vector $\mu \in \mathbb{R}^D$ of the expectation values for each dimension.

Σ The *covariance matrix* containing the variance for each dimension on its main axis and the covariances on the other spaces. It is symmetric and defined as:

$$\Sigma = \begin{bmatrix} \text{Var}(x_1) & \text{Cov}(x_1, x_2) & \cdots & \text{Cov}(x_1, x_D) \\ \text{Cov}(x_2, x_1) & \text{Var}(x_2) & \cdots & \text{Cov}(x_2, x_D) \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}(x_D, x_1) & \text{Cov}(x_D, x_2) & \cdots & \text{Var}(x_D) \end{bmatrix} \quad (3.27)$$

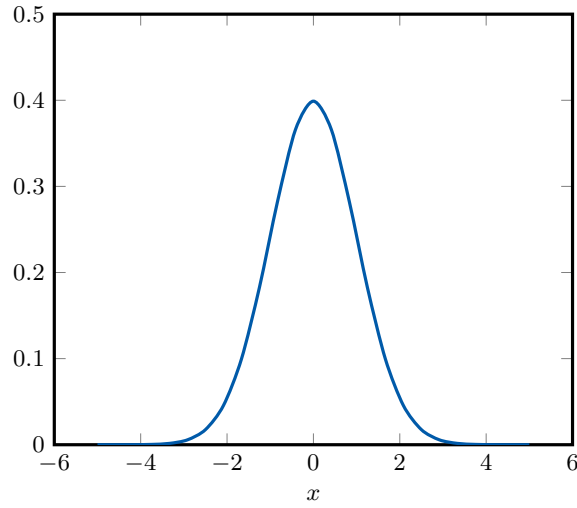


Figure 3.4.: Standard Gaussian Distribution $\mathcal{N}(x | 0, 1)$

Properties

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \Sigma) = \frac{1}{\sqrt{2\pi}^D \sqrt{\det \Sigma}} \exp \left\{ -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}) \right\} \quad (3.28)$$

$$\mathbb{E}(\mathbf{x}) = \boldsymbol{\mu} \quad (3.29)$$

$$\text{Var}(x_i) = \Sigma_{ii} \quad (3.30)$$

Geometry

Moments

3.2.5. Partitioned Gaussian Distributions

3.3. Central Limit Theorem

The distribution of the sum of N i.i.d. random variables becomes increasingly Gaussian as N increases. That is, with $N \rightarrow \infty$, it converges towards a Gaussian.

3.4. Probability Rules

Joint Distribution

$$p(x, y) \quad (3.31)$$

Marginal Distribution

$$p(y) = \int p(x, y) \, dx \quad (3.32)$$

Conditional Distribution

$$p(y | x) = \frac{p(x, y)}{p(x)} \quad (3.33)$$

Probabilistic/Stochastic Independence

$$p(x, y) = p(x)p(y) \quad (3.34)$$

Chain Rule of Probabilities

$$p(x_1, \dots, x_n) = p(x_1 | x_2, \dots, x_n) p(x_2, \dots, x_n) \quad (3.35)$$

$$= p(x_1 | x_2, \dots, x_n) p(x_2 | x_3, \dots, x_n) \cdots p(x_{n-1} | x_n) p(x_n) \quad (3.36)$$

Bayes Rule

$$p(y | x) = \frac{p(x | y) p(y)}{p(x)} \quad (3.37)$$

- Posterior: $p(y | x)$
- Likelihood: $p(x | y)$
- Prior: $p(y)$
- Normalization Factor: $p(x) = \int p(x, y) dy = \int p(x | y) p(y) dy$

3.5. Expectation, Variance and Moments

3.5.1. Expectation

The *expectation* value of a random variable x with a distribution $p(x)$ is defined as:

$$\mathbb{E}_{x \sim p(x)}(f(x)) = \mathbb{E}_x(f(x)) = \mathbb{E}(f(x)) = \begin{cases} \sum_x f(x) p(x) & \text{for discrete distributions} \\ \int_x f(x) p(x) dx & \text{for continuous distributions} \end{cases} \quad (3.38)$$

This gives a similar formula for the *conditional expectation*:

$$\mathbb{E}_{x \sim p(x | y)}(f(x)) = \mathbb{E}_x(f(x)) = \mathbb{E}(f(x)) = \begin{cases} \sum_x f(x) p(x | y) & \text{for discrete distributions} \\ \int_x f(x) p(x | y) dx & \text{for continuous distributions} \end{cases} \quad (3.39)$$

With enough samples, the expectation value can be approximated using the arithmetic mean:

$$\mathbb{E}(f(x)) \approx \frac{1}{N} \sum_{i=1}^N f(x_i) \quad (3.40)$$

Calculation Rules Let x, y be random variables and $\alpha \in \mathbb{R}$.

$$\mathbb{E}(\alpha x) = \alpha \mathbb{E}(x) \quad (3.41)$$

$$\mathbb{E}(x + y) = \mathbb{E}(x) + \mathbb{E}(y) \quad (3.42)$$

$$\mathbb{E}(xy) = \mathbb{E}(x) \mathbb{E}(y) \quad (3.43)$$

Equation 3.43 only holds if x and y are statistically independent.

3.5.2. Variance and Covariance

The *variance* measures the spread of the variable in relation to its mean:

$$\text{Var}(x) = \mathbb{E}((x - \mathbb{E}(x))^2) = \mathbb{E}(x^2) - (\mathbb{E}(x))^2 \quad (3.44)$$

The *covariance* measures the correlation between two variables (how much the variables change together):

$$\text{Cov}(x, y) = \mathbb{E}_{x,y}(xy) - \mathbb{E}_x(x) \mathbb{E}_y(y) \quad (3.45)$$

$$\text{Cov}(\mathbf{x}, \mathbf{y}) = \mathbb{E}_{\mathbf{x}, \mathbf{y}}(\mathbf{x} \mathbf{y}^T) - \mathbb{E}_{\mathbf{x}}(\mathbf{x}) \mathbb{E}_{\mathbf{y}}(\mathbf{y}^T) \quad (3.46)$$

This gives the following very important rule (with $\boldsymbol{\mu}$ and Σ from the Gaussian):

$$\mathbb{E}(\mathbf{x} \mathbf{x}^T) = \boldsymbol{\mu} \boldsymbol{\mu}^T + \Sigma \quad (3.47)$$

3.5.3. Moments

A *moment* is defined as

$$m_n = E(x^n) \quad (3.48)$$

The *central moment* is defined as

$$cm_n = E((x - \mu)^n) \quad (3.49)$$

which leads to another definition of the variance, skewness and kurtosis:

cm_2 Variance (measure of spreading)

cm_3 Skewness (measure of asymmetry)

cm_4 Kurtosis (measure of heavy/light tailed-ness)

3.6. Exponential Family

The *exponential family* is a large class of distributions that are analytically interesting, because taking the log of them simplifies them a lot. All distributions of this family are unimodal with the following general form:

$$p(\mathbf{x} | \boldsymbol{\eta}) = h(\mathbf{x})g(\boldsymbol{\eta}) \exp \{ \boldsymbol{\eta}^T \mathbf{u}(\mathbf{x}) \} \quad (3.50)$$

where $\boldsymbol{\eta}$ is the natural parameter and

$$g(\boldsymbol{\eta}) \int_{-\infty}^{+\infty} h(\mathbf{x}) \exp \{ \boldsymbol{\eta}^T \mathbf{u}(\mathbf{x}) \} = 1 \quad (3.51)$$

holds. g can be interpreted as a normalization to make this property hold true.

3.6.1. Example: Bernoulli Distribution

The Bernoulli distribution is part of the exponential family and decomposes as

$$\text{Bern}(x | \mu) = \mu^x (1 - \mu)^{1-x} \quad (3.52)$$

$$= \exp \{ x \ln(\mu) + (1 - x) \ln(1 - \mu) \} \quad (3.53)$$

$$= (1 - \mu) \exp \left\{ \ln \left(\frac{\mu}{1 - \mu} \right) x \right\} \quad (3.54)$$

with the logistic sigmoid

$$\sigma(\eta) = \frac{1}{1 + \exp(-\eta)} \quad (3.55)$$

and

$$\eta = \ln \left(\frac{\mu}{1 - \mu} \right) \quad (3.56)$$

we can write the Bernoulli distribution as

$$p(x | \mu) = \sigma(-\eta) \exp(\eta x) \quad (3.57)$$

which, in the exponential family form, gives:

$$h(x) = 1 \quad g(\eta) = \sigma(-\eta) \quad u(x) = x \quad (3.58)$$

3.6.2. Example: Gaussian Distribution

The Gaussian distribution is part of the exponential family and decomposes as

$$\mathcal{N}(x | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{1}{2\sigma^2} (x - \mu)^2 \right\} \quad (3.59)$$

$$= \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{1}{2\sigma^2} x^2 + \frac{\mu}{\sigma^2} x - \frac{\mu^2}{2\sigma^2} \right\} \quad (3.60)$$

$$= h(x)g(\boldsymbol{\eta}) \exp \{ \boldsymbol{\eta}^T \mathbf{u}(x) \} \quad (3.61)$$

with

$$\boldsymbol{\eta} = \left[-\frac{1}{2\sigma^2} \quad \frac{\mu}{\sigma^2} \right]^T \quad h(x) = 1 \quad g(\boldsymbol{\eta}) = \sqrt{-\frac{\eta_1}{\pi}} \exp \left\{ \frac{\eta_2^2}{4\eta_1} \right\} \quad \mathbf{u}(x) = \begin{bmatrix} x^2 \\ x \end{bmatrix} \quad (3.62)$$

3.7. Information Theory and Entropy

Information theory is about how to represent information compactly (as few bits as possible) and therefore about compression. This raises three questions:

- How to measure complexity?
- How to measure the “distance” between probability distributions?
- How to reconstruct data?

3.7.1. Information and Entropy

- Information is hiding in data.
- E.g. in the English alphabet, every letter has a different probability p_i of occurring.
- A lower probability indicates that the data point contains more information.
- The *average information*, called *entropy* can be calculated as

$$H(p) = - \sum_i p_i \log_2(p_i) \quad (3.63)$$

3.7.2. Kullback-Leibler Divergence

The *Kullback-Leibler Divergence* is a similarity measurement between probability distributions, defined by

$$\text{KL}(p \parallel q) = - \int p(x) \ln(q(x)) \, dx - \left(- \int p(x) \ln(p(x)) \, dx \right) \quad (3.64)$$

$$= - \int p(x) \ln \left(\frac{q(x)}{p(x)} \right) \, dx \quad (3.65)$$

The KL divergence represents the average additional number of bits required to specify a symbol x , if the underlying probability distribution is the estimated $q(x)$ and not the true one $p(x)$.

Some properties:

- $\text{KL}(p \parallel q) \neq \text{KL}(q \parallel p)$ not a distance
- $\text{KL}(p \parallel q) \geq 0$ non-negative distance
- $(\forall x : p(x) = q(x)) \implies \text{KL}(p \parallel q) = 0$

There exist other metrics than KL, but KL is deeply connected to maximum likelihood estimation.

3.8. Wrap-Up

- Random variables (both continuous and discrete)
- Probability distributions
- Basic rules of probability theory
- Expectation and variance
- Gaussian distribution and its importance
- Information and entropy

4. Fundamentals: Optimization

“All learning problems are essentially optimization problems on data” (Christopher G. Atkeson, Professor at CMU)

All machine learning problems are optimization problems in the form

$$\min_{\theta} J(\theta, \mathcal{D}) \quad (4.1)$$

$$\text{s.t.} \quad f(\theta, \mathcal{D}) = 0 \quad (4.2)$$

$$g(\theta, \mathcal{D}) \geq 0 \quad (4.3)$$

with parameters θ to enable learning, a data set \mathcal{D} to learn from, a cost function $J(\theta, \mathcal{D})$ to measure the performance and equality and inequality constraints $f(\theta, \mathcal{D}) = 0$, $g(\theta, \mathcal{D}) \geq 0$.

4.1. Convexity

A set $C \subseteq \mathbb{R}^n$ is *convex* iff for all $\mathbf{x}, \mathbf{y} \in C$ and for all $\alpha \in [0, 1]$ the following holds:

$$\alpha \mathbf{x} + (1 - \alpha) \mathbf{y} \in C \quad (4.4)$$

Intuition: Every point on a line drawn between two arbitrary points in space lie in the set itself. The set has no “bays”.

A function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is *convex* iff for all $\mathbf{x}, \mathbf{y} \in \text{Domain}(f)$ and for all $\alpha \in [0, 1]$ the following holds:

$$f(\alpha \mathbf{x} + (1 - \alpha) \mathbf{y}) \leq \alpha f(\mathbf{x}) + (1 - \alpha) f(\mathbf{y}) \quad (4.5)$$

Intuition: The drawn line between two arbitrary points on the function do not cross the function (they may touch, so linear functions are also convex).

If f is differentiable, it is convex iff for all $\mathbf{x}, \mathbf{y} \in \text{Domain}(f)$ the following holds:

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla_{\mathbf{x}} f(\mathbf{x}) (\mathbf{y} - \mathbf{x}) \quad (4.6)$$

If f is twice differentiable, it is convex iff for all $\mathbf{x} \in \text{Domain}(f)$ the following holds:

$$\nabla_{\mathbf{x}}^2 f(\mathbf{x}) \succeq 0 \quad (4.7)$$

Warning: Differentiability is not a condition for convexity!

4.2. Cost Functions

An ideal cost function is convex. But most of the time, they are not. . .

4.2.1. Common Cost Functions

Table 4.1 lists common cost functions for classification, regression, density estimation and clustering.

4.3. Constrained/Unconstrained Optimization

The general form of a constrained optimization problem is

$$\max_{\theta} J(\theta) \quad (4.8)$$

$$\text{s.t.} \quad \mathbf{f}(\theta) = 0 \quad (4.9)$$

$$\mathbf{g}(\theta) \geq 0 \quad (4.10)$$

with a cost function $J(\theta)$, some equality constraints $\mathbf{f}(\theta)$ and inequality constraints $\mathbf{g}(\theta)$.

| Problem | Example Cost Functions | Resulting Method |
|--------------------|--|-------------------------------|
| Classification | $\min_{\theta} \sum_{i=1}^n \ln \left(1 + \exp \left(- y_i x_i^T \theta \right) \right)$ | Logistic Regression |
| | $\min_{\theta_1, \theta_2} \sum_{i=1}^n \left(y_i - g \left(\theta_2^T g \left(\theta_1^T x_i \right) \right) \right)^2$ | Neural Network Classification |
| Regression | $\min_{\theta} \ \theta\ ^2 + C \sum_{i=1}^n \xi_i$ s.t. $\xi_i - (1 - y_i x_i^T \theta) \geq 0, \xi_i \geq 0$ | Support Vector Machines |
| | $\min_{\theta} \sum_{i=1}^n (y_i - \phi(x_i)^T \theta)^2$ | Linear Regression |
| | $\min_{\theta_1, \theta_2, \theta_3} \sum_{i=1}^n \left(y_i - \theta_3^T g \left(\theta_2^T g \left(\theta_1^T x_i \right) \right) \right)^2$ | Neural Network Regression |
| Density Estimation | $\min_{\theta} \sum_{i=1}^n \ln (p(x_i \theta))$ | General Formulation |
| Clustering | $\min_{\mu_1, \dots, \mu_k} \sum_{j=1}^k \sum_{i \in C_j} \ x_i - \mu_i\ ^2$ | |

Table 4.1.: Common Cost Functions

4.4. Lagrange Multipliers

With a constrained optimization problem in the general form, the *Lagrangian* is defined as

$$\mathcal{L}(\theta, \lambda, \mu, \epsilon) = J(\theta) + \lambda^T \mathbf{f}(\theta) + \mu^T (\mathbf{g}(\theta) + \epsilon^2) \quad (4.11)$$

The coefficients λ and μ are called *Lagrangian Multipliers*, the variables ϵ are called *slack variables* and are used to convert the inequality constraints into equality constraints.

To solve the optimization problem, take the derivatives w.r.t. θ , λ and μ and set them to zero:

$$\nabla_{\theta} \mathcal{L} = 0 \quad \nabla_{\lambda} \mathcal{L} = 0 \quad \nabla_{\mu} \mathcal{L} = 0 \quad (4.12)$$

If this results in any $\epsilon_i = 0$, the inequality constraint is called *active* and the solution lies on the edge of that constraint. To check whether the result really is a minima/maxima, take the second derivative of the cost function w.r.t. θ , $\nabla_{\theta}^2 J(\theta)$, and check whether the resulting Hessian is positive or negative definite, yielding that the found solution is a minima or maxima, respectively.

4.4.1. Dual Formulation

Given the so-called *primal problem*

$$\min_{\theta} J(\theta) \quad (4.13)$$

$$\text{s.t.} \quad \mathbf{f}(\theta) = 0 \quad (4.14)$$

$$\mathbf{g}(\theta) \geq 0 \quad (4.15)$$

with the Lagrangian

$$\mathcal{L}(\theta, \lambda, \mu, \epsilon) = J(\theta) + \lambda^T \mathbf{f}(\theta) + \mu^T (\mathbf{g}(\theta) + \epsilon^2) \quad (4.16)$$

the *dual problem* is

$$\max_{\lambda, \mu} \hat{\mathcal{L}}(\lambda, \mu, \epsilon) = \min_{\theta} \mathcal{L}(\theta, \lambda, \mu, \epsilon^2) \quad (4.17)$$

$$\text{s.t.} \quad \lambda \geq 0 \quad (4.18)$$

$$\mu \geq 0 \quad (4.19)$$

- If λ^* is the solution for the dual problem, then $\hat{\mathcal{L}}(\lambda^*)$ is a *lower bound* for the primal problem due to two concepts:
 - *Minimax inequality*: For any function with two arguments $\phi(x, y)$, the maximin is less or equal to the minimax:

$$\max_y \min_x \phi(x, y) \leq \min_x \max_y \phi(x, y) \quad (4.20)$$

- *Weak duality*: The primal values are always greater or equal to the dual values:

$$\min_{\theta} \max_{\substack{\lambda \geq 0 \\ \mu \geq 0}} \mathcal{L}(\theta, \lambda, \mu) \geq \max_{\substack{\lambda \geq 0 \\ \mu \geq 0}} \min_{\theta} \mathcal{L}(\theta, \lambda, \mu) \quad (4.21)$$

- In machine learning, the dual is often far more useful than the primal.
- That is because $\hat{\mathcal{L}}$ is a concave function and easy to optimize, even if J and the constraints may be nonconvex.
- Given some λ and μ , the dual is an unconstrained problem.

4.4.2. Example

Given the following optimization problem (in the real numbers):

$$\arg \max_{x,y} J(x,y) = x + y \quad (4.22)$$

$$\text{s.t. } x^2 + y^2 - 1 = 0 \quad (4.23)$$

$$2 - x \geq 0 \quad (4.24)$$

the Lagrangian is written as

$$\mathcal{L}(x, y, \lambda, \mu, \epsilon) = x + y + \lambda(x^2 + y^2 - 1) + \mu(2 - x + \epsilon^2) \quad (4.25)$$

Take the derivatives:

$$\nabla_x \mathcal{L} = 1 + 2\lambda x - \mu \quad (4.26)$$

$$\nabla_y \mathcal{L} = 1 + 2\lambda y \quad (4.27)$$

$$\nabla_\lambda \mathcal{L} = x^2 + y^2 - 1 \quad (4.28)$$

$$\nabla_\mu \mathcal{L} = 2 - x + \epsilon^2 \quad (4.29)$$

$$\nabla_\epsilon \mathcal{L} = 2\mu\epsilon \quad (4.30)$$

Settings them to zero gives the insight that either $\mu = 0$ or $\epsilon = 0$ must be true. This must be done per case.

Case 1: $\mu = 0$ This yields the following equation system:

$$0 = 1 + 2\lambda x \quad (4.31)$$

$$0 = 1 + 2\lambda y \quad (4.32)$$

$$0 = x^2 + y^2 - 1 \quad (4.33)$$

$$0 = 2 - x + \epsilon^2 \quad (4.34)$$

with the solution $x = y = \pm \frac{1}{\sqrt{2}}$ and $\epsilon^2 = \frac{1}{\sqrt{2}} - 2$, so the inequality constraint is not active as the solution fulfills the equation $x < 2$.

Case 2: $\epsilon = 0$ This yields the following equation system:

$$0 = 1 + 2\lambda x - \mu \quad (4.35)$$

$$0 = 1 + 2\lambda y \quad (4.36)$$

$$0 = x^2 + y^2 - 1 \quad (4.37)$$

$$0 = 2 - x \quad (4.38)$$

Which yields the following solutions:

$$x_1 = 2 \quad x_2 = 2 \quad (4.39)$$

$$y_1 = -i\sqrt{3} \quad y_2 = i\sqrt{3} \quad (4.40)$$

As $\epsilon = 0$, the solution must lie on the edge of the inequality constraint, thus $x = 2$. As only real solutions were wanted, this solution can be discarded.

4.5. Numerical Optimization

For a lot of optimization problems, the solution cannot be computed analytically, so these have to be approximated using numerical optimization.

The performance of numerical methods can be measured with the following questions:

- Does the algorithm converge to the optimal solution?
- How many steps does it take to converge?
- Is the convergence smooth or bumpy?
- Does it work for all types of functions or just a special type (e.g. convex)?