

## 1. Statistical Learning

### 1.1. Definition Statistical Model

Statistical Model:	$\{\mathbb{X}, \mathbb{F}, P_\theta; \theta \in \Theta\}$
Sample Space:	$\Omega$
Observation Space:	$\mathbb{X}$
Sigma Algebra:	$\mathbb{F}$
Probability:	$P_\theta$
Test:	$T: \mathbb{X} \mapsto \{\theta_0, \theta_1\}, x \mapsto T(x)$
Null Hypothesis:	$H_0: \theta \in \Theta_0$
Alternative Hypothesis:	$H_1: \theta \in \Theta_1$

**Cost Criterion  $G_T$ :**  
 $G_T: \{\theta_0, \theta_1\} \mapsto [0, 1], \theta \mapsto P(\{T(X) = 1\}|\theta)$   
 $= E[T(X); \theta] = \int_{\mathbb{X}} T(x) f_X(x|\theta) dx$

**Error Level  $\alpha$ :**  $G_T(\theta_0) \leq \alpha$

**Two Error Types:**  
False Alarm:  $\theta = \theta_0, T(x) = 1$   
 $G_T(\theta_0) = P(\{T(X) = 1\}|\theta_0)$   
Detection Error:  $\theta = \theta_1, T(x) = 0$   
 $1 - G_T(\theta_1) = P(\{T(X) = 0\}|\theta_1)$

### 1.2. Maximum Likelihood Test

**ML Ratio Test Statistic:**  

$$R(x) = \begin{cases} \frac{f_X(x|\theta_1)}{f_X(x|\theta_0)} & ; f_X(x|\theta_0) > 0 \\ \infty & ; f_X(x|\theta_0) = 0 \text{ and } f_X(x|\theta_1) > 0 \end{cases}$$

**ML Test:**  
 $T_{ML}: \mathbb{X} \mapsto \{0, 1\}, x \mapsto \begin{cases} 1 & ; R(x) > c \\ 0 & ; \text{otherwise} \end{cases}$

if  $c \neq 1$  False Alarm Error Probability can be adjusted  $\rightarrow$  Neyman Pearson Test

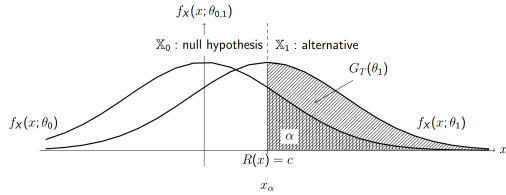
### 1.3. Neyman-Pearson-Test

The best test of  $P_0$  against  $P_1$  is

$$T_{NP}(x) = \begin{cases} 1 & R(x) > c \\ \gamma & R(x) = c \\ 0 & R(x) < c \end{cases} \quad \text{Likelihood-Ratio: } R(x) = \frac{f_X(x|\theta_1)}{f_X(x|\theta_0)}$$

$\gamma = \frac{\alpha - P_0(\{R > c\})}{P_0(\{R = c\})}$  Errorlevel  $\alpha$

Steps: For  $\alpha$  calculate  $x_\alpha$ , then  $c = R(x_\alpha)$



**Maximum Likelihood Detector:**  $T_{ML}(x) = \begin{cases} 1 & R(x) > 1 \\ 0 & \text{otherwise} \end{cases}$

**ROC Graphs:** plot  $G_T(\theta_1)$  as a function of  $G_T(\theta_0)$

### 1.4. Bayes Test (MAP Test)

Prior knowledge about possible hypotheses:

$$P(\{\theta \in \Theta_0\}) + P(\{\theta \in \Theta_1\}) = 1$$

$$T_{\text{Bayes}} = \underset{T}{\operatorname{argmin}} \{P_e\} = \begin{cases} 1 & ; \frac{f_X(x|\theta_1)}{f_X(x|\theta_0)} > c \\ 0 & ; \text{otherwise} \end{cases}$$

$$= \begin{cases} 1 & ; P(\theta_1|x) > P(\theta_0|x) \\ 0 & ; \text{otherwise} \end{cases}$$

with:

$$P_e = P(\theta_0)G_T(\theta_0) + P(\theta_1)(1 - G_T(\theta_1)), \quad c = \frac{P(\theta_0)}{P(\theta_1)}$$

if  $P(\theta_0) = P(\theta_1) \rightarrow T_{\text{Bayes}} = T_{ML}$

**Multiple Hypothesis**  $\{\theta_0, \dots, \theta_k\}; \mathbb{X}_0, \dots, \mathbb{X}_k \in \mathbb{X}$ :

$$T_{\text{Bayes}} = \underset{k \in 1, \dots, K}{\operatorname{argmin}} \{P(\theta_k|x)\}$$

**Loss Function:**

$$L(T(x), \theta) = \begin{cases} L_0 & ; T(x) = 1, \text{ but } \theta = \theta_0 \quad (\text{FALSE ALARM}) \\ L_1 & ; T(x) = 0, \text{ but } \theta = \theta_1 \quad (\text{DETEC. ERROR}) \\ 0 & ; \text{otherwise} \end{cases}$$

$L_i$  denotes the Loss Value in cases where the correct decision parameter  $\theta_i$  is missed.

$$\text{Risk}(T) = E[L(T(X), \theta)] = E[E[L(T(x), \theta)|x = X]]$$

### 1.5. Linear Alternative Tests

Estimate normal vector  $\underline{w}^\top$  and  $w_0$ , which separate  $\mathbb{X}$  into  $\mathbb{X}_0$  and  $\mathbb{X}_1$

$$\log R(\underline{x}) = -\frac{1}{2} \ln \left( \frac{\det(\underline{C}_1)}{\det(\underline{C}_0)} \right) - \frac{1}{2} (\underline{x} - \underline{\mu}_1)^\top \underline{C}_1^{-1} (\underline{x} - \underline{\mu}_1) + \frac{1}{2} (\underline{x} - \underline{\mu}_0)^\top \underline{C}_0^{-1} (\underline{x} - \underline{\mu}_0) = \ln \left( \frac{P(\theta \in \Theta_0)}{P(\theta \in \Theta_1)} \right) \quad (\text{separating surface})$$

For Gaussian  $f_X(x; \mu_k, C_k)$  with  $\theta_0$  and  $\theta_1$  corresponding to  $\{\mu_0, C_0\}$  and  $\{\mu_1, C_1\}$ , it follows that

- if  $C_0 \neq C_1$ ,  $\log R(x) = 0$  is non-linear and the separating surfaces are surfaces of second order: parabolic, hyperbolic, or elliptic surfaces.
- if  $C_0 = C_1$ ,  $\log R(x) = 0$  is affine and thus defines a hyperplane in  $\mathbb{X}$  which decomposes  $\mathbb{X}$  into  $\mathbb{X}_0$  and  $\mathbb{X}_1$ , i.e.,  

$$T: \mathbb{X} \rightarrow \mathbb{R}, \underline{x} \mapsto \begin{cases} 1 & \underline{w}^\top \underline{x} > w_0 \\ 0 & \text{otherwise} \end{cases}$$

– case 1:  $\underline{C}_0 = \underline{C}_1 = \sigma^2 \underline{I}_N$

$$\underline{w}^\top = (\underline{\mu}_1 - \underline{\mu}_0)^\top,$$

$$w_0 = \frac{1}{2} (\underline{\mu}_1^\top \underline{\mu}_1 - \underline{\mu}_0^\top \underline{\mu}_0) - \sigma^2 \ln \left( \frac{P(\theta \in \Theta_1)}{P(\theta \in \Theta_0)} \right)$$

$\underline{w}$  colinear with  $(\underline{\mu}_1 - \underline{\mu}_0)$

$\rightarrow$  hyperplane orthogonal to  $(\underline{\mu}_1 - \underline{\mu}_0)$

– case 2:  $\underline{C}_0 = \underline{C}_1 = \underline{C}$

$$\underline{w}^\top = (\underline{\mu}_1 - \underline{\mu}_0)^\top \underline{C}^{-1},$$

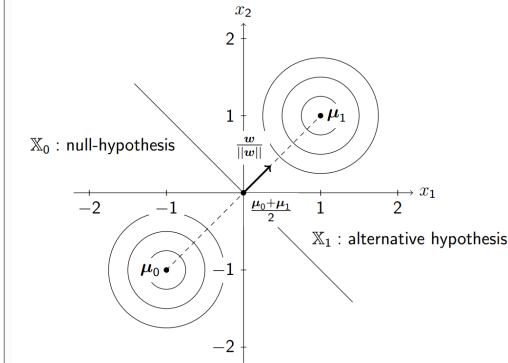
$$w_0 = \frac{1}{2} (\underline{\mu}_1 - \underline{\mu}_0)^\top \underline{C}^{-1} (\underline{\mu}_1 + \underline{\mu}_0) - \ln \left( \frac{P(\theta \in \Theta_1)}{P(\theta \in \Theta_0)} \right)$$

in general  $\underline{w}$  not colinear with  $(\underline{\mu}_1 - \underline{\mu}_0)$

$\rightarrow$  hyperplane not orthogonal to  $(\underline{\mu}_1 - \underline{\mu}_0)$

- if  $C_0 = C_1$  and  $\mu_0 = -\mu_1$ ,  $\log R(x) = 0$  is linear and defines a separating hyperplane in  $\mathbb{X}$  which contains the origin, i.e.,

$$T: \mathbb{X} \rightarrow \mathbb{R}, \underline{x} \mapsto \begin{cases} 1 & \underline{w}^\top \underline{x} > 0 \\ 0 & \text{otherwise} \end{cases}$$



## 2. Learning and Generalization

### 2.1. Empirical Risk Function and Generalization Error

ML scenarios (unknown Stochastic Model) base learning on:  
 $Risk_{emp}(T; \mathbb{S}) = \frac{1}{M} \sum_{i=1}^M L(T(\underline{x}_i), y_i), \quad (\underline{x}_i, y_i) \in \mathbb{S}$

$$\underline{x} \mapsto T(\underline{x}; \mathbb{S}) \quad T = \underset{T' \in \mathbb{T}}{\operatorname{argmin}} \{Risk_{emp}(T'; \mathbb{S})\}$$

**good Generalization:**  $Risk_{emp}(T; \mathbb{S}_{test})$  similar to  $Risk_{emp}(T; \mathbb{S})$   
**bad Generalization:**

- small  $\mathbb{T}$  that does not cover  $T_{opt} \rightarrow$  cannot be selected by ML  
 $\Rightarrow$  strong mismatch between the desired and derived Test and refers to a sort of **Bias Error Term**
- too rich  $\mathbb{T} \rightarrow$  fluctuating of the available data (measurement noise) is interpreted as meaningful information  
 $\Rightarrow$  **Overfitting**; leads to an increased **Variance Error Term**

### 2.2. Bias-Variance Decomposition

$$Risk = E_{\mathbb{S}, X, Y} [L(T(X; \mathbb{S}), Y)] = E_X [1 - P_{Y|X}(Y = T_B(X))] + (1 - P_{\mathbb{S}|X}(T(X; \mathbb{S}) = T_B(X))) (2P_{Y|X}(Y = T_B(X)) - 1),$$

$T_B(X)$  is the unknown Bayes Test  
If the potential set  $\mathbb{S}$  would be selected from a distribution such that the derived Test  $T(\underline{x}; \mathbb{S})$  and the corresponding Bayes Test  $T_B(\underline{x})$  are identical almost surely, then the Risk Function achieves its minimum value which is equal to the **Irreducible Error**  $E_X [1 - P_{Y|X}(Y = T_B(X))]$  (denotes the probability that for a given input  $\underline{x}$  the Bayes Test  $T_B(X)$  decides for the false label  $y$ ).

## 3. Classification Trees and Random Forests

### 3.1. CART Algorithms

### 3.2. Random Forests

## 4. Hypothesis Testing

making a decision based on the observations

### 4.1. Definition

Null hypothesis  $H_0: \theta \in \Theta_0$  (Assumed first to be true)

Alternate hypothesis  $H_1: \theta \in \Theta_1$  (The one to proof)

Decision rule  $\varphi: \mathbb{X} \rightarrow [0, 1]$  with

$\varphi(x) = 1$ : decide for  $H_1$ ,  $\varphi(x) = 0$ : decide for  $H_0$  Error level  $\alpha$  with  $E[d(X)|\theta] \leq \alpha, \forall \theta \in \Theta_0$

Error Type	Decision \ Reality	$H_1$ false ( $H_0$ true)	$H_1$ true ( $H_0$ false)
1 (FA) False Alarm	$H_1$ rejected ( $H_0$ accepted)	True Negative $P = 1 - \alpha$	False Negative (Type 2) $P = \beta$
2 (DE) Detection Error	$H_1$ accepted ( $H_0$ rejected)	False Positive (Type 1) $P = \alpha$	True Positive $P = 1 - \beta$

Power: Sensitivity/Recall/Hit Rate:  $\frac{TP}{TP+FN} = 1 - \beta$

Specificity/True negative rate:  $\frac{TN}{FP+TN} = 1 - \alpha$

Precision/Positive Prediction rate:  $\frac{TP}{TP+FP}$

Accuracy:  $\frac{TP+TN}{P+N} = \frac{2-\alpha-\beta}{2}$

#### 4.1.1. Design of a test

Cost criterion  $G_\varphi: \Theta \rightarrow [0, 1], \theta \mapsto E[d(X)|\theta]$

False Positive lower than  $\alpha$ :  $G_d(\theta)|_{\theta \in \Theta_0} \leq \alpha, \forall \theta \in \Theta_0$

False Negative small as possible:  $\max\{G_d(\theta)|_{\theta \in \Theta_1}\}, \forall \theta \in \Theta_1$

### 4.2. Sufficient Statistics

Sufficiency for a test  $T(X)$  means that no other test statistic, i.e., function of the observations  $\underline{x}$ , contains additional information about the parameter  $\theta$  to be estimated:

$$f_{X|T}(x|T(x) = t, \theta) = f_{X|T}(x|T(x) = t)$$

## 5. Math

$$\pi \approx 3.14159 \quad e \approx 2.71828 \quad \sqrt{2} \approx 1.414 \quad \sqrt{3} \approx 1.732$$

**Binome, Trinome**

$$(a \pm b)^2 = a^2 \pm 2ab + b^2 \quad a^2 - b^2 = (a - b)(a + b)$$

$$(a \pm b)^3 = a^3 \pm 3a^2b + 3ab^2 \pm b^3$$

$$(a + b + c)^2 = a^2 + b^2 + c^2 + 2ab + 2ac + 2bc$$

**Folgen und Reihen**

$$\sum_{k=1}^n k = \frac{n(n+1)}{2}$$

Aritmetische Summenformel

$$\sum_{k=0}^n q^k = \frac{1-q^{n+1}}{1-q}$$

Geometrische Summenformel

$$\sum_{n=0}^{\infty} \frac{z^n}{n!} = e^z$$

Exponentialreihe

**Mittelwerte** ( $\sum$  von  $i$  bis  $N$ ) (Median: Mitte einer geordneten Liste)

$$\bar{x}_{ar} = \frac{1}{N} \sum x_i \geq \bar{x}_{geo} = \sqrt[N]{\prod x_i} \geq \bar{x}_{hm} = \frac{N}{\sum \frac{1}{x_i}}$$

Arithmetisches Mittel Geometrisches Mittel Harmonisches Mittel

**Ungleichungen:** Bernoulli-Ungleichung:  $(1+x)^n \geq 1+nx$

$$||x| - |y|| \leq |x \pm y| \leq |x| + |y| \quad \left| \underline{x}^T \cdot \underline{y} \right| \leq \|\underline{x}\| \cdot \|\underline{y}\|$$

Dreiecksungleichung Cauchy-Schwarz-Ungleichung

**Mengen:** De Morgan:  $\overline{A \cap B} = \overline{A} \cup \overline{B} \quad \overline{A \cup B} = \overline{A} \cap \overline{B}$

**5.1. Exp. und Log.**  $e^x := \lim_{n \rightarrow \infty} \left(1 + \frac{x}{n}\right)^n \quad e \approx 2,71828$

$$a^x = e^{x \ln a} \quad \log_a x = \frac{\ln x}{\ln a} \quad \ln x \leq x - 1$$

$$\ln(x^a) = a \ln(x) \quad \ln\left(\frac{x}{a}\right) = \ln x - \ln a \quad \log(1) = 0$$

**5.2. Matrizen**  $\underline{A} \in \mathbb{K}^{m \times n}$

$\underline{A} = (a_{ij}) \in \mathbb{K}^{m \times n}$  hat  $m$  Zeilen (Index  $i$ ) und  $n$  Spalten (Index  $j$ )

$$(\underline{A} + \underline{B})^T = \underline{A}^T + \underline{B}^T \quad (\underline{A} \cdot \underline{B})^T = \underline{B}^T \cdot \underline{A}^T$$

$$(\underline{A}^T)^{-1} = (\underline{A}^{-1})^T \quad (\underline{A} \cdot \underline{B})^{-1} = \underline{B}^{-1} \cdot \underline{A}^{-1}$$

$$\dim \mathbb{K} = n = \text{rang } \underline{A} + \dim \ker \underline{A} \quad \text{rang } \underline{A} = \text{rang } \underline{A}^T$$

**5.2.1. Quadratische Matrizen**  $\underline{A} \in \mathbb{K}^{n \times n}$

regulär/invertierbar/nicht-singulär  $\Leftrightarrow \det(\underline{A}) \neq 0 \Leftrightarrow \text{rang } \underline{A} = n$

singulär/nicht-invertierbar  $\Leftrightarrow \det(\underline{A}) = 0 \Leftrightarrow \text{rang } \underline{A} \neq n$

orthogonal  $\Leftrightarrow \underline{A}^T = \underline{A}^{-1} \Rightarrow \det(\underline{A}) = \pm 1$

symmetrisch:  $\underline{A} = \underline{A}^T$  schiefssymmetrisch:  $\underline{A} = -\underline{A}^T$

**5.2.2. Determinante von  $\underline{A} \in \mathbb{K}^{n \times n}$ :**  $\det(\underline{A}) = |\underline{A}|$

$$\det \begin{bmatrix} \underline{A} & \underline{0} \\ \underline{C} & \underline{D} \end{bmatrix} = \det \begin{bmatrix} \underline{A} & \underline{B} \\ \underline{0} & \underline{D} \end{bmatrix} = \det(\underline{A}) \det(\underline{D})$$

$$\det(\underline{A}) = \det(\underline{A}^T) \quad \det(\underline{A}^{-1}) = \det(\underline{A})^{-1}$$

$$\det(\underline{A}\underline{B}) = \det(\underline{A}) \det(\underline{B}) = \det(\underline{B}) \det(\underline{A}) = \det(\underline{B}\underline{A})$$

Hat  $\underline{A}$  2 linear abhäng. Zeilen/Spalten  $\Rightarrow |\underline{A}| = 0$

**5.2.3. Eigenwerte (EW)  $\lambda$  und Eigenvektoren (EV)  $\underline{v}$**

$\underline{A}\underline{v} = \lambda \underline{v} \quad \det \underline{A} = \prod \lambda_i \quad \text{Sp } \underline{A} = \sum a_{ii} = \sum \lambda_i$

Eigenwerte:  $\det(\underline{A} - \lambda \underline{1}) = 0$  Eigenvektoren:  $\ker(\underline{A} - \lambda_i \underline{1}) = \underline{v}_i$

**5.2.4. Spezialfall  $2 \times 2$  Matrix  $\underline{A}$**

$$\det(\underline{A}) = ad - bc \quad \begin{bmatrix} a & b \\ c & d \end{bmatrix}^{-1} = \frac{1}{\det \underline{A}} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}$$

$$\text{Sp}(\underline{A}) = a + d$$

$$\lambda_{1/2} = \frac{\text{Sp } \underline{A}}{2} \pm \sqrt{\left(\frac{\text{Sp } \underline{A}}{2}\right)^2 - \det \underline{A}}$$

**5.2.5. Differentiation**

$$\frac{\partial \underline{x}^T \underline{y}}{\partial \underline{x}} = \frac{\partial \underline{y}^T \underline{a}}{\partial \underline{x}} = \underline{y} \quad \frac{\partial \underline{x}^T \underline{A} \underline{x}}{\partial \underline{x}} = (\underline{A} + \underline{A}^T) \underline{x}$$

$$\frac{\partial \underline{x}^T \underline{A} \underline{y}}{\partial \underline{A}} = \underline{x} \underline{y}^T \quad \frac{\partial \det(\underline{B} \underline{A} \underline{C})}{\partial \underline{A}} = \det(\underline{B} \underline{A} \underline{C}) \left( \underline{A}^{-1} \right)^T$$

**5.2.6. Ableitungsregeln** ( $\forall \lambda, \mu \in \mathbb{R}$ )

Linearität:  $(\lambda f + \mu g)'(x) = \lambda f'(x) + \mu g'(x_0)$

Produkt:  $(f \cdot g)'(x) = f'(x)g(x) + f(x)g'(x)$

Quotient:  $\left(\frac{f}{g}\right)'(x) = \frac{g(x)f'(x) - f(x)g'(x)}{g(x)^2} \quad \left(\frac{\text{NAZ-ZAN}}{N^2}\right)$

Kettenregel  $(f(g(x)))' = f'(g(x))g'(x)$

**5.3. Integrale**  $\int e^x dx = e^x = (e^x)'$

Partielle Integration:  $\int u w' = u w - \int u' w$

Substitution:  $\int f(g(x))g'(x) dx = \int f(t) dt$

$F(x) - C$	$f(x)$	$f'(x)$
$\frac{1}{q+1} x^{q+1}$	$x^q$	$q x^{q-1}$
$\frac{2\sqrt{ax^3}}{3}$	$\sqrt{ax}$	$\frac{a}{2\sqrt{ax}}$
$x \ln(ax) - x$	$\ln(ax)$	$\frac{1}{x}$
$\frac{1}{a^2} e^{ax} (ax - 1)$	$x \cdot e^{ax}$	$e^{ax} (ax + 1)$
$\frac{e^x}{\ln(a)}$	$a^x$	$a^x \ln(a)$
$-\cos(x)$	$\sin(x)$	$\cos(x)$
$\cosh(x)$	$\sinh(x)$	$\cosh(x)$
$-\ln \cos(x) $	$\tan(x)$	$\frac{1}{\cos^2(x)}$

$$\int e^{at} \sin(bt) dt = e^{at} \frac{a \sin(bt) + b \cos(bt)}{a^2 + b^2}$$

$$\int \frac{dt}{\sqrt{at+b}} = \frac{2\sqrt{at+b}}{a} \quad \int t^2 e^{at} dt = \frac{(ax-1)^2+1}{a^3} e^{at}$$

$$\int t e^{at} dt = \frac{at-1}{a^2} e^{at} \quad \int x e^{ax^2} dx = \frac{1}{2a} e^{ax^2}$$

**5.3.1. Volumen und Oberfläche von Rotationskörpern um x-Achse**

$$V = \pi \int_a^b f(x)^2 dx \quad O = 2\pi \int_a^b f(x) \sqrt{1 + f'(x)^2} dx$$

## 6. Probability Theory Basics

**6.1. Kombinatorik**

Mögliche Variationen/Kombinationen um  $k$  Elemente von maximal  $n$  Elementen zu wählen bzw.  $k$  Elemente auf  $n$  Felder zu verteilen:

	Mit Reihenfolge	Reihenfolge egal
Mit Wiederholung	$n^k$	$\binom{n+k-1}{k}$
Ohne Wiederholung	$\frac{n!}{(n-k)!}$	$\binom{n}{k}$

Permutation von  $n$  mit jeweils  $k$  gleichen Elementen:  $\frac{n!}{k_1! \cdot k_2! \cdot \dots}$

Binomialkoeffizient  $\binom{n}{k} = \binom{n}{n-k} = \frac{n!}{k! \cdot (n-k)!}$

$$\binom{n}{0} = 1 \quad \binom{n}{1} = n \quad \binom{n}{2} = 6 \quad \binom{n}{5} = 10 \quad \binom{n}{6} = 15$$

**6.2. Der Wahrscheinlichkeitsraum  $(\Omega, \mathbb{F}, P)$**

Ergebnismenge	$\Omega = \{\omega_1, \omega_2, \dots\}$	Ergebnis $\omega_j \in \Omega$
Ereignisalgebra	$\mathbb{F} = \{A_1, A_2, \dots\}$	Ereignis $A_i \subseteq \Omega$
Wahrscheinlichkeitsmaß	$P : \mathbb{F} \rightarrow [0, 1]$	$P(A) = \frac{ A }{ \Omega }$

**6.3. Wahrscheinlichkeitsmaß P**

$$P(A) = \frac{|A|}{|\Omega|} \quad P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

**6.3.1. Axiome von Kolmogorow**

Nichtnegativität:  $P(A) \geq 0 \Rightarrow P : \mathbb{F} \mapsto [0, 1]$

Normiertheit:  $P(\Omega) = 1$

Additivität:  $P\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} P(A_i)$ , wenn  $A_i \cap A_j = \emptyset, \forall i \neq j$

**6.4. Bedingte Wahrscheinlichkeit**

Bedingte Wahrscheinlichkeit für  $A$  falls  $B$  bereits eingetreten ist:

$$P_B(A) = P(A|B) = \frac{P(A \cap B)}{P(B)}$$

**6.4.1. Totale Wahrscheinlichkeit und Satz von Bayes**

Es muss gelten:  $\bigcup_{i \in I} B_i = \Omega$  für  $B_i \cap B_j = \emptyset, \forall i \neq j$

Totale Wahrscheinlichkeit:  $P(A) = \sum_{i \in I} P(A|B_i) P(B_i)$

Satz von Bayes:  $P(B_k|A) = \frac{P(A|B_k) P(B_k)}{\sum_{i \in I} P(A|B_i) P(B_i)}$

**Multiplikationssatz:**  $P(A \cap B) = P(A|B) P(B) = P(B|A) P(A)$

**6.5. Zufallsvariable**

$X : \Omega \mapsto \mathbb{R}'$  ist Zufallsvariable, wenn für jedes Ereignis  $A' \in \mathbb{F}'$  im Bildraum ein Ereignis  $A$  im Urbildraum  $\mathbb{F}$  existiert, sodass  $\{\omega \in \Omega | X(\omega) \in A'\} \in \mathbb{F}$

**6.6. Distribution**

Bezeichnung	Abk.	Zusammenhang
Wahrscheinlichkeitsdichte	pdf	$f_X(x) = \frac{dF_X(x)}{dx}$
Kumulative Verteilungsfkt.	cdf	$F_X(x) = \int_{-\infty}^x f_X(\xi) d\xi$

Joint CDF:  $F_{X,Y}(x, y) = P(\{X \leq x, Y \leq y\})$

**6.7. Relations zwischen  $f_X(x), f_{X,Y}(x, y), f_{X|Y}(x|y)$**

$$\underbrace{\int_{-\infty}^{\infty} f_{X,Y}(x, y) d\xi}_{\text{Marginalization}} = \underbrace{\int_{-\infty}^{\infty} f_{X|Y}(x, \xi) f_Y(\xi) d\xi}_{\text{Total Probability}} = f_X(x)$$

Joint PDF

**6.8. Bedingte Zufallsvariablen**

Ereignis A gegeben:  $F_{X|A}(x|A) = P(\{X \leq x\} | A)$

ZV Y gegeben:  $F_{X|Y}(x|y) = P(\{X \leq x\} | \{Y = y\})$

$$p_{X|Y}(x|y) = \frac{p_{X,Y}(x, y)}{p_Y(y)}$$

$$f_{X|Y}(x|y) = \frac{f_{X,Y}(x, y)}{f_Y(y)} = \frac{dF_{X|Y}(x|y)}{dx}$$

**6.9. Unabhängigkeit von Zufallsvariablen**

$X_1, \dots, X_n$  sind stochastisch unabhängig, wenn für jedes  $\underline{x} \in \mathbb{R}^n$  gilt:

$$F_{X_1, \dots, X_n}(x_1, \dots, x_n) = \prod_{i=1}^n F_{X_i}(x_i)$$

$$p_{X_1, \dots, X_n}(x_1, \dots, x_n) = \prod_{i=1}^n p_{X_i}(x_i)$$

$$f_{X_1, \dots, X_n}(x_1, \dots, x_n) = \prod_{i=1}^n f_{X_i}(x_i)$$

## 7. Common Distributions

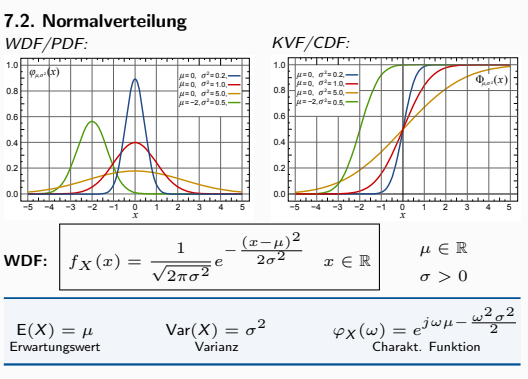
**7.1. Binomialverteilung  $B(n, p)$  mit  $p \in [0, 1], n \in \mathbb{N}$**

Folge von  $n$  Bernoulli-Experimenten

$p$ : Wahrscheinlichkeit für Erfolg  $k$ : Anzahl der Erfolge

$$p_X(k) = B_{n,p}(k) = \begin{cases} \binom{n}{k} p^k (1-p)^{n-k} & k \in \{0, \dots, n\} \\ 0 & \text{sonst} \end{cases}$$

$E[X] = np$ Erwartungswert	$\text{Var}[X] = np(1-p)$ Varianz	$G_X(z) = (pz + 1 - p)^n$ Wahrscheinlichkeitserz. Funktion
-------------------------------	--------------------------------------	---



**7.3. Sonstiges**

**Gammadistribution**  $\Gamma(\alpha, \beta): E[X] = \frac{\alpha}{\beta}$

**Exponential:**  $f(x, \lambda) = \lambda e^{-\lambda x} \quad E[X] = \lambda^{-1} \quad \text{Var}[X] = \lambda^{-2}$

## 8. Wichtige Parameter

**8.1. Erwartungswert (1. zentrales Moment)**

gibt den mittleren Wert einer Zufallsvariablen an

$$\mu_X = E[X] = \sum_{x \in \Omega'} x \cdot P_X(x) \stackrel{\triangle}{=} \int_{\mathbb{R}} x \cdot f_X(x) dx$$

diskrete  $X: \Omega \rightarrow \Omega'$       stetige  $X: \Omega \rightarrow \mathbb{R}$

$E[\alpha X + \beta Y] = \alpha E[X] + \beta E[Y] \quad X \leq Y \Rightarrow E[X] \leq E[Y]$

$E[X^2] = \text{Var}[X] + E[X]^2$

$E[X Y] = E[X] E[Y]$ , falls  $X$  und  $Y$  stochastisch unabhängig

Umkehrung nicht möglich: Unkorreliertheit  $\nRightarrow$  Stoch. Unabhängig!

**8.1.1. Für Funktionen von Zufallsvariablen  $g(x)$**

$$E[g(X)] = \sum_{x \in \Omega'} g(x) P_X(x) \stackrel{\triangle}{=} \int_{\mathbb{R}} g(x) f_X(x) dx$$

**8.2. Varianz (2. zentrales Moment)**

ist ein Maß für die Stärke der Abweichung vom Erwartungswert

$$\sigma_X^2 = \text{Var}[X] = E[(X - E[X])^2] = E[X^2] - E[X]^2$$

$\text{Var}[\alpha X + \beta] = \alpha^2 \text{Var}[X] \quad \text{Var}[X] = \text{Cov}[X, X]$

$\text{Var}\left[\sum_{i=1}^n X_i\right] = \sum_{i=1}^n \text{Var}[X_i] + \sum_{j \neq i} \text{Cov}[X_i, X_j]$

**Standard Abweichung:**  $\sigma = \sqrt{\text{Var}[X]}$

**8.3. Kovarianz**

Maß für den linearen Zusammenhang zweier Variablen

$$\begin{aligned} \text{Cov}[X, Y] &= E[(X - E[X])(Y - E[Y])^T] = \\ &= E[X Y^T] - E[X] E[Y]^T = \text{Cov}[Y, X] \end{aligned}$$

$\text{Cov}[\alpha X + \beta, \gamma Y + \delta] = \alpha \gamma \text{Cov}[X, Y]$

$\text{Cov}[X + U, Y + V] = \text{Cov}[X, Y] + \text{Cov}[X, V] + \text{Cov}[U, Y] + \text{Cov}[U, V]$

**8.3.1. Korrelation = standardisierte Kovarianz**

$$\rho(X, Y) = \frac{\text{Cov}[X, Y]}{\sqrt{\text{Var}[X] \cdot \text{Var}[Y]}} = \frac{C_{x,y}}{\sigma_x \cdot \sigma_y} \quad \rho(X, Y) \in [-1; 1]$$

**8.3.2. Kovarianzmatrix für  $\underline{z} = (\underline{x}, \underline{y})^T$**

$$\text{Cov}[\underline{z}] = \underline{C}_{\underline{z}} = \begin{bmatrix} C_X & C_{XY} \\ C_{XY} & C_Y \end{bmatrix} = \begin{bmatrix} \text{Cov}[X, X] & \text{Cov}[X, Y] \\ \text{Cov}[Y, X] & \text{Cov}[Y, Y] \end{bmatrix}$$

Immer symmetrisch:  $C_{xy} = C_{yx}!$  Für Matrizen:  $\underline{C}_{\underline{x}\underline{y}} = \underline{C}_{\underline{y}\underline{x}}^T$

9. Estimation

9.1. Estimation

Statistic Estimation treats the problem of inferring underlying characteristics of unknown random variables on the basis of observations of outputs of those random variables.

Sample Space $\Omega$	nonempty set of outputs of experiment
Sigma Algebra $\mathbb{F} \subseteq 2^\Omega$	set of subsets of outputs (events)
Probability $P : \mathbb{F} \mapsto [0, 1]$	
Random Variable $X : \Omega \mapsto \mathbb{X}$	mapped subsets of $\Omega$
Observations: $x_1, \dots, x_N$	single values of $X$
Observation Space $\mathbb{X}$	possible observations of $X$
Unknown parameter $\theta \in \Theta$	parameter of propability function
Estimator $\bigcirc \bullet \bullet : \mathbb{X} \mapsto \Theta$	$\bigcirc \bullet \bullet (X) = \hat{\theta}$ , finds $\hat{\theta}$ from $X$

unknown parm. $\theta$	estimation of param. $\hat{\theta}$
R.V. of param. $\Theta$	estim. of R.V. of parm $T(X) = \hat{\Theta}$

9.2. Quality Properties of Estimators

Consistent:  $\lim_{N \rightarrow \infty} \bigcirc \bullet \bullet (x_1, \dots, x_N) = \theta$   
Bias Bias( $\bigcirc \bullet \bullet$ ) :=  $E[\bigcirc \bullet \bullet (X_1, \dots, X_N)] - \theta$   
unbiased if Bias( $\bigcirc \bullet \bullet$ ) = 0 (biased estimators can provide better estimates than unbiased estimators.)  
Variance  $\text{Var}[\bigcirc \bullet \bullet] := E[(\bigcirc \bullet \bullet - E[\bigcirc \bullet \bullet])^2]$

9.3. Mean Square Error (MSE)

The MSE is an extension of the Variance  $\text{Var}[\bigcirc \bullet \bullet] := E[(\bigcirc \bullet \bullet - E[\bigcirc \bullet \bullet])^2]$ :

$$\varepsilon[\bigcirc \bullet \bullet] = E[(\bigcirc \bullet \bullet - \theta)^2] \stackrel{\text{MSE}}{=} \text{Var}(\bigcirc \bullet \bullet) + (\text{Bias}[\bigcirc \bullet \bullet])^2 = E[(\hat{\theta} - \theta)^2]$$

If  $\Theta$  is also r.v.  $\Rightarrow$  mean over both (e.g. Bayes est.):  
Mean MSE:  $E[(\bigcirc \bullet \bullet (X) - \Theta)^2] = E[E[(\bigcirc \bullet \bullet (X) - \Theta)^2 | \Theta = \theta]]$

9.3.1. Minimum Mean Square Error (MMSE)

Minimizes mean square error:  $\arg \min_{\hat{\theta}} E[(\hat{\theta} - \theta)^2]$   
 $E[(\hat{\theta} - \theta)^2] = E[\theta^2] - 2\hat{\theta} E[\theta] + \hat{\theta}^2$   
Solution:  $\frac{d}{d\hat{\theta}} E[(\hat{\theta} - \theta)^2] \stackrel{!}{=} 0 = -2 E[\theta] + 2\hat{\theta} \Rightarrow \hat{\theta}_{\text{MMSE}} = E[\theta]$

9.4. Maximum Likelihood

Given model  $\{\mathbb{X}, \mathbb{F}, P_\theta; \theta \in \Theta\}$ , assume  $P_\theta(\underline{x})$  or  $f_X(\underline{x}, \theta)$  for observed data  $\underline{x}$ . Estimate parameter  $\theta$  so that the likelihood  $L(\underline{x}, \theta)$  or  $L(\theta | X = \underline{x})$  to obtain  $\underline{x}$  is maximized.

**Likelihood Function:** (Prob. for  $\theta$  given  $\underline{x}$ )

Discrete:  $L(x_1, \dots, x_N; \theta) = P_\theta(x_1, \dots, x_N)$   
Continuous:  $L(x_1, \dots, x_N; \theta) = f_{X_1, \dots, X_N}(x_1, \dots, x_N, \theta)$   
If  $N$  observations are Identically Independently Distributed (i.i.d.):

$$L(\underline{x}, \theta) = \prod_{i=1}^N P_\theta(x_i) = \prod_{i=1}^N f_{X_i}(x_i)$$

**ML Estimator** (Picks  $\theta$ ):  $\bigcirc \bullet \bullet_{\text{ML}} : X \mapsto \arg \max_{\theta \in \Theta} \{L(X, \theta)\} = \arg \max_{\theta \in \Theta} \{\log L(X, \theta)\} \stackrel{\text{i.i.d.}}{=} \arg \max_{\theta \in \Theta} \{\sum \log L(x_i, \theta)\}$

Find Maximum:  $\frac{\partial L(\underline{x}, \theta)}{\partial \theta} = \frac{d}{d\theta} \log L(x; \theta) \Big|_{\theta=\hat{\theta}} \stackrel{!}{=} 0$   
Solve for  $\theta$  to obtain ML estimator function  $\hat{\theta}_{\text{ML}}$

Check quality of estimator with MSE  
Maximum-Likelihood Estimator is Asymptotically Efficient. However, there might be not enough samples and the likelihood function is often not known.

9.5. Uniformly Minimum Variance Unbiased (UMVU) Estimators (Best unbiased estimators)

Best unbiased estimator: Lowest Variance of all estimators.  
Fisher's Information Inequality: Estimate lower bound of variance if

- $L(x, \theta) > 0, \forall x, \theta$
- $L(x, \theta)$  is diffable for  $\theta$
- $\int_{\mathbb{X}} \frac{\partial}{\partial \theta} L(x, \theta) dx = \frac{\partial}{\partial \theta} \int_{\mathbb{X}} L(x, \theta) dx$

**Score Function:**  
 $g(x, \theta) = \frac{\partial}{\partial \theta} \log L(x, \theta) = \frac{\frac{\partial}{\partial \theta} L(x, \theta)}{L(x, \theta)} \quad E[g(x, \theta)] = 0$

**Fischer Information:**  
 $I_F(\theta) := \text{Var}[g(X, \theta)] = E[g(x, \theta)^2] = -E\left[\frac{\partial^2}{\partial \theta^2} \log L(X, \theta)\right]$

**Cramér-Rao Lower Bound (CRB):** (if  $\bigcirc \bullet \bullet$  is unbiased)

$$\text{Var}[\bigcirc \bullet \bullet (X)] \geq \left(\frac{\partial E[\bigcirc \bullet \bullet (X)]}{\partial \theta}\right)^2 \frac{1}{I_F(\theta)}$$
$$\text{Var}[\bigcirc \bullet \bullet (X)] \geq \frac{1}{I_F(\theta)}$$

For  $N$  i.i.d. observations:  $I_F^{(N)}(x, \theta) = N \cdot I_F^{(1)}(x, \theta)$

9.5.1. Exponential Models

If  $f_X(x) = \frac{h(x) \exp(a(\theta)t(x))}{\exp(b(\theta))}$  then  $I_F(\theta) = \frac{\partial a(\theta)}{\partial \theta} \frac{\partial E[t(X)]}{\partial \theta}$

**Some Derivations:** (check in exam)  
Uniformly: Not diffable  $\Rightarrow$  no  $I_F(\theta)$   
Normal  $\mathcal{N}(\theta, \sigma^2)$ :  $g(x, \theta) = \frac{(x-\theta)}{\sigma^2} \quad I_F(\theta) = \frac{1}{\sigma^2}$   
Binomial  $\mathcal{B}(\theta, K)$ :  $g(x, \theta) = \frac{x}{\theta} - \frac{K-x}{1-\theta} \quad I_F(\theta) = \frac{K}{\theta(1-\theta)}$

9.6. Bayes Estimation (Conditional Mean)

A Priori information about  $\theta$  is known as probability  $f_\Theta(\theta; \sigma)$  with random variable  $\Theta$  and parameter  $\sigma$ . Now the conditional pdf  $f_{X|\Theta}(x, \theta)$  is used to find  $\theta$  by minimizing the mean MSE instead of uniformly MSE.  
Mean MSE for  $\Theta$ :  $E[E[(T(X) - \Theta)^2 | \Theta = \theta]]$

**Conditional Mean Estimator:**  
 $T_{\text{CM}} : x \mapsto E[\Theta | X = x] = \int_{\Theta} \theta \cdot f_{\Theta|X}(\theta|x) d\theta$   
Posterior  $f_{\Theta|X}(\underline{x}|\underline{\theta}) = \frac{f_{X|\Theta}(\underline{x})f_\Theta(\theta)}{\int_{\Theta} f_{X|\xi}(\underline{x}, \xi) d\xi} = \frac{f_{X|\theta}(\underline{x})f_\theta(\theta)}{f_X(x)}$

**Hint:** to calculate  $f_{\Theta|X}(\theta|\underline{x})$ : Replace every factor not containing  $\theta$ , such as  $\frac{1}{f_X(x)}$  with a factor  $\gamma$  and determine  $\gamma$  at the end such that  $\int_{\Theta} f_{\Theta|X}(\theta|\underline{x}) d\theta = 1$   
MMSE:  $E[\text{Var}[X | \Theta = \theta]]$

**Multivariate Gaussian:**  $X, \Theta \sim \mathcal{N} \Rightarrow \sigma_X^2 = \sigma_{X|\Theta=\theta}^2 + \sigma_\Theta$   
 $\bigcirc \bullet \bullet_{\text{CM}} : x \mapsto E[\Theta | X = x] = \underline{\mu}_\Theta + \mathcal{C}_{\Theta, X} \mathcal{C}_X^{-1} (\underline{x} - \underline{\mu}_X)$   
MMSE:  
 $E[\|\bigcirc \bullet \bullet_{\text{CM}} - \Theta\|_2^2] = \text{tr}(\mathcal{C}_{\Theta|X}) = \text{tr}(\mathcal{C}_\Theta - \mathcal{C}_{\Theta, X} \mathcal{C}_X^{-1} \mathcal{C}_{X, \Theta})$

**Orthogonality Principle:**  
 $\bigcirc \bullet \bullet_{\text{CM}}(\underline{X}) - \Theta \perp h(\underline{X}) \Rightarrow E[(T_{\text{CM}}(\underline{X}) - \Theta)h(\underline{X})] = 0$   
**MMSE Estimator:**  $\hat{\theta}_{\text{MMSE}} = \arg \min_{\theta \in \Theta} \text{MSE}$   
minimizes the MSE for all estimators

9.7. Example:

Estimate mean  $\theta$  of  $X$  with prior knowledge  $\theta \in \Theta \sim \mathcal{N}$ :  
 $X \sim \mathcal{N}(\theta, \sigma_X^2 |_{\Theta=\theta})$  and  $\Theta \sim \mathcal{N}(m, \sigma_\Theta^2)$   
 $\hat{\theta}_{\text{CM}} = E[\Theta | X = \underline{x}] = \frac{N\sigma_\Theta^2}{\sigma_X^2 |_{\Theta=\theta} + N\sigma_\Theta^2} \hat{\theta}_{\text{ML}} + \frac{\sigma_X^2 |_{\Theta=\theta}}{\sigma_X^2 |_{\Theta=\theta} + N\sigma_\Theta^2} m$   
For  $N$  independent observations  $x_i$ :  $\hat{\theta}_{\text{ML}} = \frac{1}{N} \sum x_i$   
Large  $N \Rightarrow$  ML better, small  $N \Rightarrow$  CM better

10. Linear Estimation

$t$  is now the unknown parameter  $\theta$ , we want to estimate  $y$  and  $\underline{x}$  is the input vector... review regression problem  $\underline{y} = \underline{A}\underline{x}$  (we solve for  $\underline{x}$ ), here we solve for  $\underline{t}$ , because  $\underline{x}$  is known (measured)! Confusing...  
1. Training  $\rightarrow$  2. Estimation  
Training: We observe  $y$  and  $\underline{x}$  (knowing both) and then based on that we try to estimate  $y$  given  $\underline{x}$  (only observe  $\underline{x}$ ) with a linear model  $\hat{y} = \underline{x}^\top \underline{t}$

$$\text{Estimation: } \hat{y} = \underline{x}^\top \underline{t} + m \quad \text{or} \quad \hat{y} = \underline{x}^\top \underline{t}$$

Given:  $N$  observations  $(y_i, \underline{x}_i)$ , unknown parameters  $\underline{t}$ , noise  $m$   
 $\underline{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix} \quad \underline{X} = \begin{bmatrix} \underline{x}_1^\top \\ \vdots \\ \underline{x}_m^\top \end{bmatrix} \quad \text{Note: } \hat{y} \neq y!!$   
Problem: Estimate  $y$  based on given (known) observations  $\underline{x}$  and unknown parameter  $\underline{t}$  with assumed linear Model:  $\hat{y} = \underline{x}^\top \underline{t}$   
Note  $y = \underline{x}^\top \underline{t} + m \rightarrow y = \underline{x}'^\top \underline{t}'$  with  $\underline{x}' = \begin{pmatrix} \underline{x} \\ 1 \end{pmatrix}, \quad t' = \begin{pmatrix} \underline{t} \\ m \end{pmatrix}$   
Sometimes in Exams:  $\hat{y} = \underline{x}^\top \underline{t} \Leftrightarrow \underline{\hat{t}} = \underline{T}^\top \underline{y}$   
estimate  $\underline{x}$  given  $\underline{y}$  and unknown  $\underline{T}$

10.1. Least Square Estimation (LSE)

Tries to minimize the square error for linear Model:  $\hat{y}_{\text{LS}} = \underline{x}^\top \underline{t}_{\text{LS}}$   
Least Square Error:  $\min_{\underline{t}} \left[ \sum_{i=1}^N (y_i - \underline{x}_i^\top \underline{t})^2 \right] = \min_{\underline{t}} \|\underline{y} - \underline{X}\underline{t}\|$   
$$\underline{t}_{\text{LS}} = (\underline{X}^\top \underline{X})^{-1} \underline{X}^\top \underline{y}$$

$\hat{y}_{\text{LS}} = \underline{X} \underline{t}_{\text{LS}} \in \text{span}(\underline{X})$   
**Orthogonality Principle:**  $N$  observations  $\underline{x}_i \in \mathbb{R}^d$   
 $\underline{Y} - \underline{X} \underline{T}_{\text{LS}} \perp \text{span}[\underline{X}] \Leftrightarrow \underline{Y} - \underline{X} \underline{T}_{\text{LS}} \in \text{null}[\underline{X}^\top]$ , thus  
 $\underline{X}^\top (\underline{Y} - \underline{X} \underline{T}_{\text{LS}}) = 0$  and if  $N > d \wedge \text{rang}[\underline{X}] = d$ :  
 $\underline{T}_{\text{LS}} = (\underline{X}^\top \underline{X})^{-1} \underline{X}^\top \underline{Y}$

10.2. Linear Minimum Mean Square Estimator (LMMSE)

Estimate  $y$  with linear estimator  $\underline{t}$ , such that  $\hat{y} = \underline{t}^\top \underline{x} + m$   
Note: the Model does not need to be linear! The estimator is linear!

$$\hat{y}_{\text{LMMSE}} = \arg \min_{\underline{t}, m} E[\|\underline{y} - (\underline{t}^\top \underline{x} + m)\|_2^2]$$

If Random joint variable  $\underline{z} = \begin{pmatrix} \underline{x} \\ y \end{pmatrix}$  with  
 $\underline{\mu}_{\underline{z}} = \begin{pmatrix} \underline{\mu}_{\underline{x}} \\ \mu_y \end{pmatrix}$  and  $\underline{C}_{\underline{z}} = \begin{bmatrix} \underline{C}_{\underline{x}} & \underline{c}_{xy} \\ \underline{c}_{yx} & \mu_y \end{bmatrix}$  then  
LMMSE Estimation of  $y$  given  $\underline{x}$  is  
 $\hat{y} = \mu_y + \underbrace{\underline{c}_{y\underline{x}} \underline{C}_{\underline{x}}^{-1} (\underline{x} - \underline{\mu}_{\underline{x}})}_{=\underline{t}^\top} = \underbrace{\underline{c}_{y\underline{x}} \underline{C}_{\underline{x}}^{-1} \underline{x}}_{=\underline{m}} + \underbrace{\mu_y - \underline{c}_{y\underline{x}} \underline{C}_{\underline{x}}^{-1} \underline{\mu}_{\underline{x}}}_{=0}$   
Minimum MSE:  $E[\|\underline{y} - (\underline{x}^\top \underline{t} + m)\|_2^2] = c_y - \underline{c}_{y\underline{x}} \underline{C}_{\underline{x}}^{-1} \underline{c}_{\underline{x}y}$   
**Hint:** First calculate  $\hat{y}$  in general and then set variables according to system equation.  
**Multivariate:**  $\underline{\hat{y}} = \underline{T}_{\text{LMMSE}}^\top \underline{x} \quad \underline{T}_{\text{LMMSE}}^\top = \underline{C}_{\underline{y}\underline{x}} \underline{C}_{\underline{x}}^{-1}$

If  $\underline{\mu}_{\underline{z}} = \underline{0}$  then  
Estimator  $\hat{y} = \underline{c}_{y, \underline{x}} \underline{C}_{\underline{x}}^{-1} \underline{x}$   
Minimum MSE:  $E[c_{y, \underline{x}}] = c_y - \underline{t}^\top \underline{c}_{\underline{x}, y}$

### 10.3. Matched Filter Estimator (MF)

For channel  $\underline{y} = \underline{h}x + \underline{v}$ , Filtered:  $\underline{\hat{t}}^\top \underline{y} = \underline{\hat{t}}^\top \underline{h}x + \underline{\hat{t}}^\top \underline{v}$

Find Filter  $\underline{\hat{t}}^\top$  that maximizes SNR =  $\frac{\|\underline{h}x\|}{\|\underline{v}\|}$

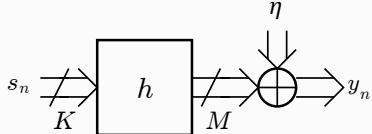
$$\underline{\hat{t}}_{MF} = \max_t \left\{ \frac{E[(\underline{\hat{t}}^\top \underline{h}x)^2]}{E[(\underline{\hat{t}}^\top \underline{v})^2]} \right\}$$

In the lecture (estimate  $\underline{h}$ ):

$$\underline{\hat{t}}_{MF} = \max_T \left\{ \frac{E[\underline{\hat{h}}^H \underline{h}]^2}{\text{tr}[\text{Var}[\underline{Tn}]]} \right\}$$

$$\underline{\hat{t}}_{MF} = \underline{T}_{MF} \underline{y} \quad \underline{T}_{MF} \propto \underline{C}_{\underline{h}} \underline{S}^H \underline{C}_{\underline{n}}^{-1}$$

### 10.4. Example



System Model:  $\underline{y}_n = \underline{H}\underline{s}_n + \eta_n$

with  $\underline{H} = (h_{m,k}) \in \mathbb{C}^{M \times K}$  ( $m \in [1, M], k \in [1, K]$ )

Linear Channel Model  $\underline{y} = \underline{S}\underline{h} + \underline{n}$  with  $\underline{h} \sim \mathcal{N}(0, \underline{C}_{\underline{h}})$  and  $\underline{n} \sim \mathcal{N}(0, \underline{C}_{\underline{n}})$

Linear Estimator  $\underline{T}$  estimates  $\underline{\hat{h}} = \underline{T}\underline{y} \in \mathbb{C}^{M \times K}$

$$\underline{T}_{MMSE} = \underline{C}_{\underline{h}} \underline{y}^{-1} = \underline{C}_{\underline{h}} \underline{S}^H (\underline{S} \underline{C}_{\underline{h}} \underline{S}^H + \underline{C}_{\underline{n}})^{-1}$$

$$\underline{T}_{ML} = \underline{T}_{Cor} = (\underline{S}^H \underline{C}_{\underline{n}}^{-1} \underline{S})^{-1} \underline{S}^H \underline{C}_{\underline{n}}^{-1}$$

$$\underline{T}_{MF} \propto \underline{C}_{\underline{h}} \underline{S}^H \underline{C}_{\underline{n}}^{-1}$$

For Assumption  $\underline{S}^H \underline{S} = N \sigma_s^2 \underline{1}_{K \times M}$  and  $\underline{C}_{\underline{n}} = \sigma_n^2 \underline{1}_{N \times M}$

Estimator	Averaged Squared Bias	Variance
ML/Correlator	0	$KM \frac{\sigma_\eta^2}{N\sigma_s^2}$
Matched Filter	$\sum_{i=1}^{KM} \lambda_i \left( \frac{\lambda_i}{\lambda_1} - 1 \right)^2$	$\sum_{i=1}^{KM} \left( \frac{\lambda_i}{\lambda_1} \right)^2 \frac{\sigma_\eta^2}{N\sigma_s^2}$
MMSE	$\sum_{i=1}^{KM} \lambda_i \left( \frac{1}{1 + \frac{\sigma_\eta^2}{\lambda_i N \sigma_s^2}} - 1 \right)^2$	$\sum_{i=1}^{KM} \frac{1}{\left( 1 + \frac{\sigma_\eta^2}{\lambda_i N \sigma_s^2} \right)^2} \frac{\sigma_\eta^2}{N\sigma_s^2}$

### 10.5. Estimators

Upper Bound: Uniform in  $[0; \theta]$ :  $\hat{\theta}_{ML} = \frac{2}{N} \sum x_i$

Probability  $p$  for  $\mathcal{B}(p, N)$ :  $\hat{p}_{ML} = \frac{\bar{x}}{N}$   $\hat{p}_{CM} = \frac{\bar{x}+1}{N+2}$

Mean  $\mu$  for  $\mathcal{N}(\mu, \sigma^2)$ :  $\hat{\mu}_{ML} = \frac{1}{N} \sum_{i=1}^N x_i$

Variance  $\sigma^2$  for  $\mathcal{N}(\mu, \sigma^2)$ :  $\hat{\sigma}_{ML}^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2$

## 11. Gaussian Stuff

### 11.1. Gaussian Channel

Channel:  $Y = h s_i + N$  with  $h \sim \mathcal{N}$ ,  $N \sim \mathcal{N}$

$$L(y_1, \dots, y_N) = \prod_{i=1}^N f_{Y_i}(y_i, h)$$

$$f_{Y_i}(y_i, h) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(y_i - h s_i)^2\right)$$

$$\hat{h}_{ML} = \underset{h}{\text{argmin}} \left\{ \left\| \underline{y} - h \underline{s} \right\|^2 \right\} = \frac{\underline{s}^\top \underline{y}}{\underline{s}^\top \underline{s}}$$

If multidimensional channel:  $\underline{y} = \underline{S}\underline{h} + \underline{n}$ :

$$L(\underline{y}, \underline{h}) = \frac{1}{\sqrt{\det(2\pi\mathbf{C})}} \exp\left(-\frac{1}{2}(\underline{y} - \underline{S}\underline{h})^\top \underline{C}^{-1}(\underline{y} - \underline{S}\underline{h})\right)$$

$$l(\underline{y}, \underline{h}) = \frac{1}{2} \left( \log(\det(2\pi\mathbf{C})) - (\underline{y} - \underline{S}\underline{h})^\top \underline{C}^{-1}(\underline{y} - \underline{S}\underline{h}) \right)$$

$$\frac{d}{dh} (\underline{y} - \underline{S}\underline{h})^\top \underline{C}^{-1}(\underline{y} - \underline{S}\underline{h}) = -2\underline{S}^\top \underline{C}^{-1}(\underline{y} - \underline{S}\underline{h})$$

Gaussian Covariance: if  $Y \sim \mathcal{N}(0, \sigma^2)$ ,  $N \sim \mathcal{N}(0, \sigma^2)$ :

$$\underline{C}_Y = \text{Cov}[Y, Y] = E[(Y - \mu)(Y - \mu)^\top] = E[Y Y^\top]$$

For Channel  $Y = Sh + N$ :  $E[Y Y^\top] = S E[h h^\top] S^\top + E[N N^\top]$

### 11.2. Multivariate Gaussian Distributions

A vector  $\underline{x}$  of  $n$  independent Gaussian random variables  $x_i$  is jointly Gaussian. If  $\underline{x} \sim \mathcal{N}(\underline{\mu}_{\underline{x}}, \underline{C}_{\underline{x}})$ :

$$f_{\underline{x}}(\underline{x}) = f_{x_1, \dots, x_n}(x_1, \dots, x_n) = \frac{1}{\sqrt{\det(2\pi\mathbf{C}_{\underline{x}})}} \exp\left(-\frac{1}{2}(\underline{x} - \underline{\mu}_{\underline{x}})^\top \underline{C}_{\underline{x}}^{-1}(\underline{x} - \underline{\mu}_{\underline{x}})\right)$$

Affine transformations  $\underline{y} = \underline{A}\underline{x} + \underline{b}$  are jointly Gaussian with

$$\underline{y} \sim \mathcal{N}(\underline{A}\underline{\mu}_{\underline{x}} + \underline{b}, \underline{A}\underline{C}_{\underline{x}}\underline{A}^\top)$$

All marginal PDFs are Gaussian as well

Contour Lines

Ellipsoid with central point  $E[\underline{y}]$  and main axis are the eigenvectors of  $\underline{C}_{\underline{y}}^{-1}$

### 11.3. Conditional Gaussian

$$\underline{A} \sim \mathcal{N}(\underline{\mu}_{\underline{A}}, \underline{C}_{\underline{A}}), \underline{B} \sim \mathcal{N}(\underline{\mu}_{\underline{B}}, \underline{C}_{\underline{B}})$$

$$\Rightarrow (A|B=b) \sim \mathcal{N}(\underline{\mu}_{A|B}, \underline{C}_{A|B})$$

Conditional Mean:

$$E[A|B=b] = \underline{\mu}_{A|B=b} = \underline{\mu}_{\underline{A}} + \underline{C}_{\underline{A}\underline{B}} \underline{C}_{\underline{B}\underline{B}}^{-1} (\underline{b} - \underline{\mu}_{\underline{B}})$$

Conditional Variance:

$$\underline{C}_{A|B} = \underline{C}_{\underline{A}\underline{A}} - \underline{C}_{\underline{A}\underline{B}} \underline{C}_{\underline{B}\underline{B}}^{-1} \underline{C}_{\underline{B}\underline{A}}$$

### 11.4. Misc

If CDF of gaussian distribution given  $\Phi(z) \sim \mathcal{N}(0, 1)$  then for  $X \sim \mathcal{N}(1, 1)$  the CDF is given as  $\Phi(x - \mu_x)$

## 12. Sequences

### 12.1. Random Sequences

Sequence of a random variable. Example: result of a dice is RV, roll a dice several times is a random sequence.

### 12.2. Markov Sequence $X_n : \Omega \rightarrow X_n$

Sequence of memoryless state transitions with certain probabilities.

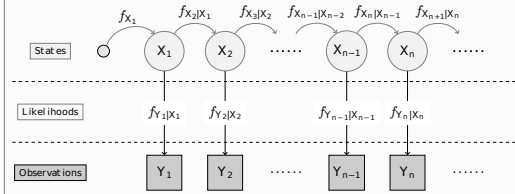
1. state:  $f_{X_1}(x_1)$

2. state:  $f_{X_2|X_1}(x_2|x_1)$

n. state:  $f_{X_n|X_{n-1}}(x_n|x_{n-1})$

### 12.3. Hidden Markov Chains

Problem: states  $X_i$  are not visible and can only be guessed indirectly as a random variable  $Y_i$ .



Conditional pdf  $f_{X_n|\underline{Y}_n}$  Likelihood pdf  $f_{Y_n|X_n}$

State-transition pdf  $f_{X_n|X_{n-1}}$

Estimation:

$$f_{X_n|\underline{Y}_n} \propto f_{Y_n|X_n} \cdot \int_{\underline{X}} f_{X_n|X_{n-1}} \cdot f_{X_{n-1}|\underline{Y}_{n-1}} d\underline{x}_{n-1}$$

## 13. Recursive Estimation

### 13.1. Kalman-Filter

recursively calculates the most likely state from previous state estimates and current observation. Shows optimum performance for Gauss-Markov Sequences.

State space:

$$\underline{x}_n = \underline{G}_n \underline{x}_{n-1} + \underline{B}_n \underline{u}_n + \underline{v}_n$$

$$\underline{y}_n = \underline{H}_n \underline{x}_n + \underline{w}_n$$

With gaussian process/measurement noise  $\underline{v}_n/\underline{w}_n$

Short notation:  $E[\underline{x}_n|\underline{y}_{n-1}] = \hat{\underline{x}}_{n|n-1}$   $E[\underline{x}_n|\underline{y}_n] = \hat{\underline{x}}_{n|n}$

$$E[\underline{y}_n|\underline{y}_{n-1}] = \hat{\underline{y}}_{n|n-1} \quad E[\underline{y}_n|\underline{y}_n] = \hat{\underline{y}}_{n|n}$$

#### 1. step: Prediction

$$\text{Mean: } \hat{\underline{x}}_{n|n-1} = \underline{G}_n \hat{\underline{x}}_{n-1|n-1}$$

$$\text{Covariance: } \underline{C}_{\underline{x}_{n|n-1}} = \underline{G}_n \underline{C}_{\underline{x}_{n-1|n-1}} \underline{G}_n^\top + \underline{C}_{\underline{v}}$$

#### 2. step: Update

$$\text{Mean: } \hat{\underline{x}}_{n|n} = \hat{\underline{x}}_{n|n-1} + \underline{K}_n (\underline{y}_n - \underline{H}_n \hat{\underline{x}}_{n|n-1})$$

$$\text{Covariance: } \underline{C}_{\underline{x}_{n|n}} = \underline{C}_{\underline{x}_{n|n-1}} + \underline{K}_n \underline{H}_n \underline{C}_{\underline{x}_{n|n-1}}$$

$$\hat{\underline{x}}_{n|n} = \hat{\underline{x}}_{n|n-1} + \underline{K}_n (\underline{y}_n - \underline{H}_n \hat{\underline{x}}_{n|n-1})$$

estimation  $E[X_n | Y_{n-1}=y_{n-1}]$  correction:  $E[X_n | \Delta Y_n=y_n]$  innovation:  $\Delta y_n$

With optimal Kalman-gain (prediction for  $\underline{x}_n$  based on  $\Delta y_n$ ):

$$\underline{K}_n = \underline{C}_{\underline{x}_{n|n-1}} \underline{H}_n^\top (\underline{H}_n \underline{C}_{\underline{x}_{n|n-1}} \underline{H}_n^\top + \underline{C}_{\underline{w}_n})^{-1}$$

$\underline{C}_{\delta y_n}$

Innovation: closeness of the estimated mean value to the real value

$$\Delta \underline{y}_n = \underline{y}_n - \hat{\underline{y}}_{n|n-1} = \underline{y}_n - \underline{H}_n \hat{\underline{x}}_{n|n-1}$$

$$\text{Init: } \hat{\underline{x}}_{0|-1} = E[X_0] \quad \sigma_{0|-1}^2 = \text{Var}[X_0]$$

$$\text{MMSE Estimator: } \hat{\underline{x}} = \int \underline{f}_n f_{X_n|Y(n)}(\underline{x}_n|\underline{y}_n) d\underline{x}_n$$

For non linear problems: Suboptimum nonlinear Filters: Extended KF, Unscented KF, ParticleFilter

### 13.2. Extended Kalman (EKF)

Linear approximation of non-linear  $g, h$

$$\underline{x}_n = g_n(\underline{x}_{n-1}, \underline{u}_n) \quad \underline{u}_n \sim \mathcal{N}$$

$$\underline{y}_n = h_n(\underline{x}_{n-1}, \underline{u}_n) \quad \underline{u}_n \sim \mathcal{N}$$

### 13.3. Unscented Kalman (UKF)

Approximation of desired PDF  $f_{X_n|Y_n}(x_n|y_n)$  by Gaussian PDF.

### 13.4. Particle-Filter

For non linear state space and non-gaussian noise

Non-linear State space:

$$\underline{x}_n = g_n(\underline{x}_{n-1}, \underline{u}_n)$$

$$\underline{y}_n = h_n(\underline{x}_{n-1}, \underline{u}_n)$$

$$\text{Posterior Conditional PDF: } f_{X_n|Y_n}(x_n|y_n) \propto f_{Y_n|X_n}(y_n|x_n) \cdot \int_{\underline{X}} \underbrace{f_{X_n|X_{n-1}}(x_n|x_{n-1})}_{\text{state transition}} \underbrace{f_{X_{n-1}|Y_{n-1}}(x_{n-1}|y_{n-1})}_{\text{last conditional PDF}} dx_{n-1}$$

$N$  random Particles with particle weight  $w_n^i$  at time  $n$

Monte-Carlo-Integration:  $I = E[g(X)] \approx I_N = \frac{1}{N} \sum_{i=1}^N \tilde{g}(x^i)$

Importance Sampling: Instead of  $f_X(x)$  use Importance Density  $q_X(x)$

$$I_N = \frac{1}{N} \sum_{i=1}^N \tilde{w}^i g(x^i) \text{ with weights } \tilde{w}^i = \frac{f_X(x^i)}{q_X(x^i)}$$

If  $\int f_{X_n}(x) dx \neq 1$  then  $I_N = \sum_{i=1}^N \tilde{w}^i g(x^i)$

### 13.5. Conditional Stochastic Independence

$$P(A \cap B|E) = P(A|E) \cdot P(B|E)$$

Given  $Y, X$  and  $Z$  are independent if

$$f_{Z|Y,X}(z|y,x) = f_{Z|Y}(z|y) \cdot f_{Z|X}(z|x)$$

$$f_{X,Z|Y}(x,z|y) = f_{Z|Y}(z|y) \cdot f_{X|Y}(x|y)$$

$$f_{Z|X,Y}(z|x,y) = f_{Z|Y}(z|y) \text{ or } f_{X|Z,Y}(x|z,y) = f_{X|Y}(x|y)$$

## 14. Hypothesis Testing

making a decision based on the observations

### 14.1. Definition

Null hypothesis  $H_0 : \theta \in \Theta_0$  (Assumed first to be true)

Alternate hypothesis  $H_1 : \theta \in \Theta_1$  (The one to proof)

Decision rule  $\varphi : \mathbb{X} \rightarrow [0, 1]$  with

$\varphi(x) = 1$ : decide for  $H_1$ ,  $\varphi(x) = 0$ : decide for  $H_0$  Error level  $\alpha$  with  $E[d(X)|\theta] \leq \alpha, \forall \theta \in \Theta_0$

Error Type	Decision \ Reality	$H_1$ false ( $H_0$ true)	$H_1$ true ( $H_0$ false)
1 (FA)	$H_1$ rejected	True Negative	False Negative (Type 2)
False Alarm	( $H_0$ accepted)	$P = 1 - \alpha$	$P = \beta$
2 (DE)	$H_1$ accepted	False Positive (Type 1)	True Positive
Detection ( $H_0$ rejected)		$P = \alpha$	$P = 1 - \beta$
Error			

$$\text{Power: Sensitivity/Recall/Hit Rate: } \frac{TP}{TP+FN} = 1 - \beta$$

$$\text{Specificity/True negative rate: } \frac{TN}{FP+TN} = 1 - \alpha$$

$$\text{Precision/Positive Prediction rate: } \frac{TP}{TP+FP}$$

$$\text{Accuracy: } \frac{TP+TN}{P+N} = \frac{2-\alpha-\beta}{2}$$

#### 14.1.1. Design of a test

Cost criterion  $G_\varphi : \Theta \rightarrow [0, 1], \theta \mapsto E[d(X)|\theta]$

False Positive lower than  $\alpha$ :  $G_d(\theta)|_{\theta \in \Theta_0} \leq \alpha, \forall \theta \in \Theta_0$

False Negative small as possible:  $\max\{G_d(\theta)|_{\theta \in \Theta_1}\}, \forall \theta \in \Theta_1$

### 14.2. Sufficient Statistics

Sufficiency for a test  $T(X)$  means that no other test statistic, i.e., function of the observations  $\underline{x}$ , contains additional information about the parameter  $\theta$  to be estimated:

$$f_{X|T}(x|T(x) = t, \theta) = f_{X|T}(x|T(x) = t)$$

## 15. Tests

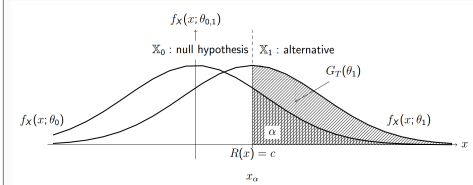
### 15.1. Neyman-Pearson-Test

The best test of  $P_0$  against  $P_1$  is

$$d_{\text{NP}}(x) = \begin{cases} 1 & R(x) > c \\ \gamma & R(x) = c \\ 0 & R(x) < c \end{cases} \quad \text{Likelihood-Ratio: } R(x) = \frac{f_X(x; \theta_1)}{f_X(x; \theta_0)}$$

$$\gamma = \frac{\alpha - P_0(\{R > c\})}{P_0(\{R = c\})} \quad \text{Errorlevel } \alpha$$

Steps: For  $\alpha$  calculate  $x_\alpha$ , then  $c = R(x_\alpha)$



$$\text{Maximum Likelihood Detector: } d_{\text{ML}}(x) = \begin{cases} 1 & R(x) > 1 \\ 0 & \text{otherwise} \end{cases}$$

ROC Graphs: plot  $G_d(\theta_1)$  as a function of  $G_d(\theta_0)$

### 15.2. Bayes Test (MAP Detector)

Prior knowledge on possible hypotheses:  $P(\{\theta \in \Theta_0\}) + P(\{\theta \in \Theta_1\}) = 1$ , minimizes the probability of a wrong decision.

$$d_{\text{Bayes}} = \begin{cases} 1 & \frac{f_X(x|\theta_1)}{f_X(x|\theta_0)} > \frac{c_0 P(\theta_0|x)}{c_1 P(\theta_1|x)} \\ 0 & \text{otherwise} \end{cases} = \begin{cases} 1 & P(\theta_1|x) > P(\theta_0|x) \\ 0 & \text{otherwise} \end{cases}$$

Risk weights  $c_0, c_1$  are 1 by default.

If  $P(\theta_0) = P(\theta_1)$ , the Bayes test is equivalent to the ML test

$$\text{Loss Function } L(d(x), \theta) = \begin{cases} c_0 & \text{type 1 } d(x) = 1, \text{ but } \theta = \theta_0 \\ c_1 & \text{type 2 } d(x) = 0, \text{ but } \theta = \theta_1 \end{cases}$$

$$\text{risk}(d) = E[L(d(X), \theta)] = E[E[L(d(x), \theta)|x = X]]$$

$$\text{Multiple Hypothesis } d_{\text{Bayes}} = \begin{cases} 0 & x \in \mathbb{X}_0 \\ 1 & x \in \mathbb{X}_1 \\ 2 & x \in \mathbb{X}_2 \end{cases}$$

### 15.3. Linear Alternative Tests

$$d : \mathbb{X} \rightarrow \mathbb{R}, \underline{x} \mapsto \begin{cases} 1 & \underline{w}^\top \underline{x} - w_0 > 0 \\ 0 & \text{otherwise} \end{cases}$$

Estimate normal vector  $\underline{w}^\top$ , which separates  $\mathbb{X}$  into  $\mathbb{X}_0$  and  $\mathbb{X}_1$

$$\log R(\underline{x}) = \frac{\ln(\det(\underline{C}_0))}{\ln(\det(\underline{C}_1))} + \frac{1}{2}(\underline{x} - \underline{\mu}_0)^\top \underline{C}_0^{-1}(\underline{x} - \underline{\mu}_0) - \frac{1}{2}(\underline{x} - \underline{\mu}_1)^\top \underline{C}_1^{-1}(\underline{x} - \underline{\mu}_1) = 0$$

For 2 Gaussians, with  $\underline{C}_0 = \underline{C}_1 = \underline{C}$ :  $\underline{w}^\top = (\underline{\mu}_1 - \underline{\mu}_0)^\top \underline{C}$

$$\text{and constant translation } w_0 = \frac{(\underline{\mu}_1 - \underline{\mu}_0)^\top \underline{C}(\underline{\mu}_1 - \underline{\mu}_0)}{2}$$

