

# MUSTERERKENNUNG

Vorlesung im Sommersemester 2017

Prof. E.G. Schukat-Talamazzini

Stand: 6. März 2017

## Teil VII

### Normalverteilungsklassifikatoren

#### Multivariate Normalverteilungsdichte

#### Normalverteilungsklassifikatoren

#### Maximum-Likelihood Parameterschätzung

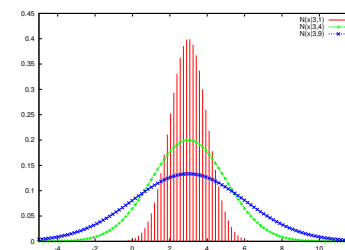
#### Maximum-a posteriori- und Bayesschätzung

#### Graphische Gaußsche Modelle

#### Mathematische Hilfsmittel

### Univariate Normalverteilungsdichte

$$\mathcal{N}(x \mid \mu, \sigma^2) \stackrel{\text{def}}{=} \frac{1}{\sigma\sqrt{2\pi}} \cdot \exp \left\{ -\frac{(x - \mu)^2}{2\sigma^2} \right\}$$



#### Definition

Eine stetige Zufallsvariable  $\mathbb{X}$  heißt (univariat) **normalverteilt** mit Mittelwert  $\mu \in \mathbb{R}$  und Varianz  $\sigma^2 \neq 0$ , wenn gilt:

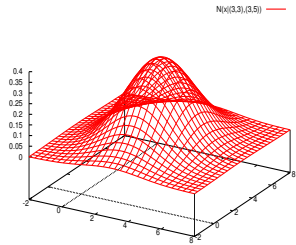
$$f_{\mathbb{X}}(x) = \mathcal{N}(x \mid \mu, \sigma^2)$$

#### Bemerkung

Unter der Annahme klassenweise *statistisch unabhängiger* und *normalverteilter* Merkmale läßt sich die (naive!) Bayesregel mit Hilfe von  $K \cdot D$  univariaten NV-Dichten realisieren.

## Bivariat unkorrelierte Normalverteilungsdichte

$$\mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}, \sigma_1^2, \sigma_2^2) \stackrel{\text{def}}{=} \frac{1}{2\pi\sigma_1\sigma_2} \cdot \exp \left\{ -\frac{1}{2} \cdot \left( \frac{(x_1 - \mu_1)^2}{\sigma_1^2} + \frac{(x_2 - \mu_2)^2}{\sigma_2^2} \right) \right\}$$



### Definition

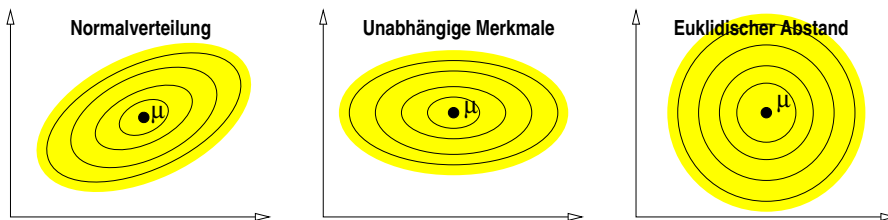
Eine stetiger Zufallsvektor  $\mathbb{X} = (\mathbb{X}_1, \mathbb{X}_2)$  heißt **bivariat unkorreliert normalverteilt** mit Mittelwertvektor  $\boldsymbol{\mu} \in \mathbb{R}^2$  und Varianzen  $\sigma_1^2, \sigma_2^2 > 0$ , wenn gilt:

$$f_{\mathbb{X}}(x_1, x_2) = \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}, \sigma_1^2, \sigma_2^2)$$

### Bemerkung

Für Normalverteilungen sind *Unkorreliertheit* und *Unabhängigkeit* äquivalent. Obige Dichte entspricht also dem Produkt  $\mathcal{N}(x_1 \mid \mu_1, \sigma_1^2) \cdot \mathcal{N}(x_2 \mid \mu_2, \sigma_2^2)$  der univariaten NV-Dichten (Randverteilungen).

## Parameterreduzierte Normalverteilungsdichten



Symmetrisch  
positiv-definit

Diagonalmatrix

Einheitsmatrix  
skaliert

$$\begin{pmatrix} \sigma_{11} & \sigma_{12} & \dots & \sigma_{1D} \\ \sigma_{21} & \sigma_{22} & \dots & \sigma_{2D} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{D1} & \sigma_{D2} & \dots & \sigma_{DD} \end{pmatrix} \quad \begin{pmatrix} \sigma_1^2 & 0 & \dots & 0 \\ 0 & \sigma_2^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_D^2 \end{pmatrix} \quad \begin{pmatrix} \sigma^2 & 0 & \dots & 0 \\ 0 & \sigma^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma^2 \end{pmatrix}$$

allgemeines  
Hyperellipsoid

Trägheitsachsen  
parallel zu  
Koordinatenachsen

skalierte  
Hypersphäre

$(D+1) \cdot \frac{D}{2}$  Parameter

$D$  Parameter

1 Parameter

## Multivariate Normalverteilungsdichte

### Definition

Ein Zufallsvektor  $\mathbb{X} = (\mathbb{X}_1, \dots, \mathbb{X}_D)^\top$  heißt **multivariat normalverteilt**, falls er der  $D$ -dimensionalen Verteilungsdichtefunktion

$$\mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}, \mathbf{S}) \stackrel{\text{def}}{=} \frac{1}{\sqrt{\det(2\pi\mathbf{S})}} \cdot \exp \left\{ -\frac{1}{2} \cdot (\mathbf{x} - \boldsymbol{\mu})^\top \mathbf{S}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}$$

gehört. Es ist  $\boldsymbol{\mu} \in \mathbb{R}^D$  der **Erwartungswertvektor** der Verteilung; die positiv-definite, symmetrische Matrix  $\mathbf{S} \in \mathbb{R}^{D \times D}$  heißt **Kovarianzmatrix** der Normalverteilung.

### Bemerkungen

1. Die Isolinen (Hyperebenen gleicher Dichtewerte) der multivariaten NV-Dichte besitzen die Form von *Hyperellipsoiden*.
2. Die Richtungen und Radien ihrer Achsen entnehmen wir den Eigenvektoren und Eigenwerten der Diagonalisierung  $\mathbf{S} = \mathbf{U}\mathbf{D}\mathbf{U}^\top$ .

## Ist $\mathcal{N}(\boldsymbol{\mu}, \mathbf{S})$ ein gutes Verteilungsmodell ?

Das kommt ganz auf die Anwendung & den Lerndatenvorrat an

### Das NV-Modell ist zu simpel für unsere Daten

- Unimodale Dichtelandschaft ? Löwe/Löwin
- Elliptische Symmetrie ? nichtnegative Merkmale
- Exponentielles Abklingverhalten ? Ausreißer

### Das NV-Modell ist zu komplex für unseren Klassifikator

- Speicheraufwand  $O(D^2 \cdot K)$  ? Bilder, Microarrays
- Rechenaufwand  $O(D^2 \cdot K)$  ? Echtzeitanwendungen
- Robustheit der Schätzung  $\hat{\mathbf{S}} = \mathbf{S}(\omega)$  ? Rang und Inversenbildung

## Multivariate Normalverteilungsdichte

### Normalverteilungsklassifikatoren

#### Maximum-Likelihood Parameterschätzung

#### Maximum-a posteriori- und Bayesschätzung

#### Graphische Gaußsche Modelle

#### Mathematische Hilfsmittel

## Normalverteilungsklassifikator

### Definition

Einen Klassifikator mit den Prüfgrößen

$$u_{\kappa}(\mathbf{x}) = P(\mathbf{x}, \Omega_{\kappa}) = p_{\kappa} \cdot \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}_{\kappa}, \mathbf{S}_{\kappa}), \quad \mathbf{x} \in \mathbb{R}^D$$

für  $\kappa = 1, \dots, K$  bezeichnet man als  $D$ -dimensionalen **Normalverteilungsklassifikator** mit den Verteilungsparametern  $[p_{\kappa}, \boldsymbol{\mu}_{\kappa}, \mathbf{S}_{\kappa}]_{\kappa=1..K}$ .

### Bemerkung

In der Praxis verwendet man einfachheitshalber die dazu *antitonen* Prüfgrößen

$$u_{\kappa}(\mathbf{x}) = -2 \cdot \log(P(\mathbf{x}, \Omega_{\kappa})),$$

die quadratische Funktionen der Mustermerkmale sind.

Entscheidungsregel: ➡ Prüfgröße **minimieren** (Minuszeichen)

## Prüfgrößen der NV-Bayesregel

Normalverteilungsklassifikator mit uneingeschränkten Kovarianzmatrizen  $\mathbf{S}_{\kappa}$

$$u_{\kappa}(\mathbf{x}) = \underbrace{-2 \log p_{\kappa} + \log |2\pi \mathbf{S}_{\kappa}|}_{\gamma_{\kappa}} + \underbrace{(\mathbf{x} - \boldsymbol{\mu}_{\kappa})^{\top} \cdot \mathbf{S}_{\kappa}^{-1} \cdot (\mathbf{x} - \boldsymbol{\mu}_{\kappa})}_{\text{Mahalanobisabstand } \|\mathbf{x} - \boldsymbol{\mu}_{\kappa}\|_{\mathbf{S}_{\kappa}}^2}$$

### Bemerkungen

1. Je Klasse  $1 + D + \binom{D+1}{2}$  Parameter ➡  $O(D^2 K)$
2. Je Muster und Klasse  $3D^2$  Addit./Multiplik. ➡  $O(D^2 K)$

$$\tilde{\mathbf{x}}^{\top} \mathbf{S}_{\kappa}^{-1} \tilde{\mathbf{x}} = \sum_{i=1}^D \sum_{j=1}^D \tilde{x}_i c_{\kappa ij} \tilde{x}_j, \quad \mathbf{C}_{\kappa} = \mathbf{S}_{\kappa}^{-1}$$

3. Für den Abstandsausdruck lohnt sich die folgende Betrachtung:

$$(\mathbf{x} - \boldsymbol{\mu}_{\kappa})^{\top} \mathbf{S}_{\kappa}^{-1} (\mathbf{x} - \boldsymbol{\mu}_{\kappa}) = \underbrace{\mathbf{x}^{\top} \mathbf{S}_{\kappa}^{-1} \mathbf{x}}_{\text{spur}(\mathbf{S}_{\kappa}^{-1} \cdot \mathbf{x} \mathbf{x}^{\top})} - \underbrace{2 \boldsymbol{\mu}_{\kappa}^{\top} \mathbf{S}_{\kappa}^{-1} \mathbf{x}}_{\mathbf{a}_{\kappa}^{\top}} + \underbrace{\boldsymbol{\mu}_{\kappa}^{\top} \mathbf{S}_{\kappa}^{-1} \boldsymbol{\mu}_{\kappa}}_{\mathbf{c}_{\kappa}}$$

## Prüfgrößen der naiven NV-Bayesregel

Normalverteilungsklassifikator mit diagonalen Kovarianzmatrizen  $\mathbf{S}_{\kappa}$

$$u_{\kappa}(\mathbf{x}) = \gamma_{\kappa} + \sum_{d=1}^D \left( \frac{x_d - \mu_{\kappa,d}}{\sigma_{\kappa,d}} \right)^2$$

mit der Konstanten

$$\gamma_{\kappa} = -2 \log p_{\kappa} + D \cdot \log(2\pi) + \sum_d \log \sigma_{\kappa,d}^2$$

### Bemerkungen

1. Je Klasse  $1 + D + D$  Parameter ➡  $O(DK)$
2. Je Muster und Klasse  $4D$  Addit./Multipl./Divis. ➡  $O(DK)$
3. Keine Merkmalkorrelationen — keine „schrägen“ Klassengebiete!

## Prüfgrößen der sphärischen NV-Bayesregel

Normalverteilungsklassifikator mit skaliertem Einheitskovarianz  $\mathbf{S}_\kappa = \sigma_\kappa^2 \mathbf{E}$

$$u_\kappa(\mathbf{x}) = \gamma_\kappa + \frac{\|\mathbf{x} - \boldsymbol{\mu}_\kappa\|^2}{\sigma_\kappa^2}$$

mit der Konstanten

$$\gamma_\kappa = -2 \log p_\kappa + D \cdot \log(2\pi) + 2D \cdot \log \sigma_\kappa$$

### Bemerkungen

1. Je Klasse  $1 + D + 1$  Parameter  $\Rightarrow O(DK)$
2. Je Muster und Klasse  $3D$  Addit./Multipl./Divis.  $\Rightarrow O(DK)$
3. Klassengebiete  $\hat{=}$  Hyperkugeln unterschiedlicher Radien

## Prüfgrößen des Minimum-Abstand-Klassifikators

Normalverteilungsklassifikator mit Einheitskovarianz  $\mathbf{S}_\kappa = \mathbf{E}$

$$u_\kappa(\mathbf{x}) = \gamma_\kappa + \|\mathbf{x} - \boldsymbol{\mu}_\kappa\|^2$$

mit der Konstanten

$$\gamma_\kappa = -2 \log p_\kappa + D \cdot \log(2\pi)$$

### Bemerkungen

1. Je Klasse  $1 + D + 0$  Parameter  $\Rightarrow O(DK)$
2. Je Muster und Klasse  $2D$  Addit./Multipl./Divis.  $\Rightarrow O(DK)$
3. Klassengebiete  $\hat{=}$  Hyperkugeln identischer Radien
4. *Modifizierter MAK* — incl. Klassengewicht  $\gamma_\kappa$
5. *Gewöhnlicher MAK* — excl. Klassengewicht  $\gamma_\kappa$

## Prüfgrößen des Mahalanobis-Klassifikators

Normalverteilungsklassifikator mit klassenunabhängiger Kovarianz  $\mathbf{S}_\kappa = \mathbf{S}_0$

$$u_\kappa(\mathbf{x}) = \gamma_\kappa + \underbrace{(\mathbf{x} - \boldsymbol{\mu}_\kappa)^\top \cdot \mathbf{S}_0^{-1} \cdot (\mathbf{x} - \boldsymbol{\mu}_\kappa)}_{\text{Mahalanobisabstand } \|\mathbf{x} - \boldsymbol{\mu}_\kappa\|_{\mathbf{S}_0}^2}$$

mit der Konstanten

$$\gamma_\kappa = -2 \log p_\kappa + D \cdot \log(2\pi) + \log |\mathbf{S}_0|$$

### Bemerkungen

1. Je Klasse  $1 + D$  Parameter zzgl.  $\mathbf{S}_0$   $\Rightarrow O(DK + D^2)$
2. Je Klasse  $2D$  Addit./Multiplik. zzgl. quadr. Form  $\Rightarrow O(DK + D^2)$
3. Für den Abstandsdruck lohnt sich die folgende Betrachtung:

$$(\mathbf{x} - \boldsymbol{\mu}_\kappa)^\top \mathbf{S}_0^{-1} (\mathbf{x} - \boldsymbol{\mu}_\kappa) = \underbrace{\mathbf{x}^\top \mathbf{S}_0^{-1} \mathbf{x}}_{\text{spur}(\mathbf{S}_0^{-1} \cdot \mathbf{x} \mathbf{x}^\top)} - \underbrace{2 \boldsymbol{\mu}_\kappa^\top \mathbf{S}_0^{-1} \mathbf{x}}_{\mathbf{a}_\kappa^\top} + \underbrace{\boldsymbol{\mu}_\kappa^\top \mathbf{S}_0^{-1} \boldsymbol{\mu}_\kappa}_{c_\kappa}$$

## Prüfgrößen des Richter-Klassifikators

Normalverteilungsklassifikator mit isotrop skaliertem Kovarianz  $\mathbf{S}_\kappa = \alpha_\kappa \mathbf{S}_0$

$$u_\kappa(\mathbf{x}) = \gamma_\kappa + \underbrace{\alpha_\kappa^{-1} \cdot (\mathbf{x} - \boldsymbol{\mu}_\kappa)^\top \cdot \mathbf{S}_0^{-1} \cdot (\mathbf{x} - \boldsymbol{\mu}_\kappa)}_{\alpha_\kappa^{-1} \cdot \|\mathbf{x} - \boldsymbol{\mu}_\kappa\|_{\mathbf{S}_0}^2}$$

mit der Konstanten

$$\gamma_\kappa = -2 \log p_\kappa + D \cdot \log(2\pi) + D \cdot \log \alpha_\kappa + \log |\mathbf{S}_0|$$

### Bemerkungen

1. Je Klasse  $1 + D + 1$  Parameter zzgl.  $\mathbf{S}_0$   $\Rightarrow O(DK + D^2)$
2. Je Klasse  $2D$  Addit./Multiplik. zzgl. quadr. Form  $\Rightarrow O(DK + D^2)$
3. Für den Abstandsdruck lohnt sich die folgende Betrachtung:

$$\mathbf{x}^\top \mathbf{S}_\kappa^{-1} \mathbf{x} = \alpha_\kappa^{-1} \cdot \underbrace{\text{spur}(\mathbf{S}_0^{-1} \cdot \mathbf{x} \mathbf{x}^\top)}_{c_x}$$

## Prüfgrößen des Eigenraumklassifikators

Normalverteilungsklassifikator mit achsenparallelen Kovarianzen  $\mathbf{S}_\kappa = \mathbf{U} \mathbf{D}_\kappa \mathbf{U}^\top$

$$u_\kappa(\mathbf{x}) = \gamma_\kappa + \underbrace{(\mathbf{U}^\top(\mathbf{x} - \boldsymbol{\mu}_\kappa))^\top \cdot \mathbf{D}_\kappa^{-1} \cdot (\mathbf{U}^\top(\mathbf{x} - \boldsymbol{\mu}_\kappa))}_{\|\mathbf{U}^\top(\mathbf{x} - \boldsymbol{\mu}_\kappa)\|_{\mathbf{D}_\kappa}^2}$$

mit der Konstanten

$$\gamma_\kappa = -2 \log p_\kappa + D \cdot \log(2\pi) + \sum_d \log \lambda_{\kappa d}$$

### Bemerkungen

1. Je Klasse  $1 + D + D$  Parameter zzgl.  $\mathbf{U}$   $\Rightarrow O(DK + D^2)$
2. Je Klasse  $4D$  Operationen für  $\|\cdot\|_{\mathbf{D}_\kappa}^2$  zzgl.  $D^2$  für  $\mathbf{U}^\top \mathbf{x}$   $\Rightarrow O(DK + D^2)$
3. Für den Abstandsausdruck lohnt sich die folgende Betrachtung:

$$\mathbf{x}^\top \mathbf{S}_\kappa^{-1} \mathbf{x} = \mathbf{x}^\top \mathbf{U} \mathbf{D}_\kappa^{-1} \mathbf{U}^\top \mathbf{x} = (\mathbf{U}^\top \mathbf{x})^\top \mathbf{D}_\kappa^{-1} (\mathbf{U}^\top \mathbf{x}) = \sum_{d=1}^D (u_d^\top \mathbf{x})^2 / \lambda_{\kappa d}$$

4. Es kommt auch eine *unvollständige* Entwicklung in Betracht, bei der Trägheitsachsen mit kleinen Eigenwerten ignoriert werden ...

## Parameterschätzung für Wahrscheinlichkeitsmodelle

Verteilungsmodell  $\Rightarrow$  Lerndaten  $\Rightarrow$  Parameterschätzwert

### Parametrische Verteilungsdichtefamilie

Die Wertetupel  $\mathbf{x} \in \mathbb{R}^D$  eines Zufallsvektors  $\mathbb{X}$  seien gemäß

$$\{f(\mathbf{x}|\boldsymbol{\theta}) \mid \boldsymbol{\theta} \in \mathcal{M}\}$$

verteilt; jede Verteilungsdichte der Familie ist durch ein Feld  $\boldsymbol{\theta}$  von Parametern aus einer Mannigfaltigkeit  $\mathcal{M}$  charakterisiert.

### Repräsentative Lernstichprobe

Die unbekannte Verteilung von  $\mathbb{X}$  ist durch eine Stichprobe  $\omega$  repräsentiert, deren Elemente  $\{\mathbf{x}_1, \dots, \mathbf{x}_T\}$  unabhängig und identisch gemäß  $f(\cdot|\boldsymbol{\theta})$  verteilt gezogen wurden.

### Problem

Wie lautet der **beste Schätzwert**  $\hat{\boldsymbol{\theta}}$  für die unbekannten Parameter  $\boldsymbol{\theta}^*$  ?

## Multivariate Normalverteilungsdichte

## Normalverteilungsklassifikatoren

## Maximum-Likelihood Parameterschätzung

## Maximum-a posteriori- und Bayesschätzung

## Graphische Gaußsche Modelle

## Mathematische Hilfsmittel

## Maximum-Likelihood Schätzung

### Lemma

Die (logarithmierte) Ziehungswahrscheinlichkeit für den unabhängigen und identischen Datensatz  $\omega$  beträgt

$$\ell_\theta(\omega) = \log \prod_{\mathbf{x} \in \omega} f(\mathbf{x}|\boldsymbol{\theta}) = \sum_{\mathbf{x} \in \omega} \log f(\mathbf{x}|\boldsymbol{\theta}) .$$

Die Größe  $\ell_\theta(\omega)$  heißt **Likelihoodfunktion** von  $\boldsymbol{\theta}$ .

### Definition

Die **Maximum-Likelihood-Schätzung** (MLS) der Parameter einer Dichtefamilie  $[f(\mathbf{x}|\boldsymbol{\theta})]$  maximiert die parameterbedingte Stichprobenwahrscheinlichkeit, d.h. es gilt

$$\hat{\boldsymbol{\theta}}_{\text{ML}} = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \prod_{\mathbf{x} \in \omega} f(\mathbf{x}|\boldsymbol{\theta}) = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \sum_{\mathbf{x} \in \omega} \log f(\mathbf{x}|\boldsymbol{\theta}) .$$

### Bemerkung

Der ML-Schätzwert  $\hat{\boldsymbol{\theta}}_{\text{ML}}$  ist von allen Parameterwerten derjenige, zu dem die vorliegenden Daten  $\omega$  am besten passen.

## Maximum-Likelihood Schätzung

### Satz

Der ML-Schätzer ist **erwartungstreu**, d.h.: ist eine Zufallsvariable  $\mathbb{X}$  gemäß  $f(\mathbf{x}|\boldsymbol{\theta}^*)$  verteilt, so ist der Erwartungswert des ML-Schätzers für eine Stichprobe unabhängiger Realisierungen von  $\mathbb{X}$  gleich  $\boldsymbol{\theta}^*$ .

### Bemerkungen

1. Für eine repräsentative Lernstichprobe zunehmenden Umfangs strebt der ML-Schätzwert gegen den wahren Parametervektor.
2. Über das Verhalten des ML-Schätzwertes bei Verwendung einer individuellen, endlichen Probe trifft der Satz keinerlei verbindliche Aussage.
3. Gehorcht der Datenerzeugungsprozeß nicht tatsächlich für irgendeinen festen Parameterwert  $\boldsymbol{\theta} \in \mathcal{M}$  dem postulierten Verteilungsgesetz  $f(\mathbf{x}|\boldsymbol{\theta})$ , so besitzen selbst die asymptotischen ML-Parameter  $\hat{\boldsymbol{\theta}}_{\text{ML}}$  keine Aussagekraft.

## ML-Schätzung für den NV-Klassifikator

### Erzeugungswahrscheinlichkeit

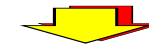
einer unabhängig und identisch verteilten, etikettierten Stichprobe

$$P\left(\bigcup_{\kappa} \omega_{\kappa}\right) = \prod_{\kappa=1}^K P(\omega_{\kappa}) = \prod_{\kappa=1}^K \prod_{\mathbf{x} \in \omega_{\kappa}} P(\Omega_{\kappa}) \cdot P(\mathbf{x}|\Omega_{\kappa})$$

### Logarithmierte ML-Zielgröße

Parametrisiert durch  $(p_{\kappa}, \boldsymbol{\theta}_{\kappa})$ ,  $\kappa = 1, \dots, K$

$$\log \prod_{\kappa=1}^K \prod_{\mathbf{x} \in \omega_{\kappa}} p_{\kappa} \cdot f(\mathbf{x}|\boldsymbol{\theta}_{\kappa}) = \sum_{\kappa=1}^K T_{\kappa} \log p_{\kappa} + \sum_{\kappa=1}^K \left( \sum_{\mathbf{x} \in \omega_{\kappa}} \log f(\mathbf{x}|\boldsymbol{\theta}_{\kappa}) \right)$$



zerfällt in  $(K + 1)$  voneinander unabhängige Optimierungsprobleme

## ML-Schätzung für den NV-Klassifikator

mit vollbesetzten klassenabhängigen Kovarianzmatrizen

### Satz

Die Maximum-Likelihood-Parameter eines Normalverteilungsklassifikators bezüglich einer etikettierten Stichprobe  $[\omega_{\kappa}]$  lauten

$$\begin{aligned} \hat{p}_{\kappa} &= T_{\kappa} / \sum_{\lambda=1}^K T_{\lambda} \\ \hat{\boldsymbol{\mu}}_{\kappa} &= \frac{1}{T_{\kappa}} \sum_{\mathbf{x} \in \omega_{\kappa}} \mathbf{x} \\ \hat{\mathbf{S}}_{\kappa} &= \frac{1}{T_{\kappa}} \sum_{\mathbf{x} \in \omega_{\kappa}} (\mathbf{x} - \hat{\boldsymbol{\mu}}_{\kappa})(\mathbf{x} - \hat{\boldsymbol{\mu}}_{\kappa})^{\top} \\ &= \frac{1}{T_{\kappa}} \sum_{\mathbf{x} \in \omega_{\kappa}} \mathbf{x} \mathbf{x}^{\top} - \hat{\boldsymbol{\mu}}_{\kappa} \hat{\boldsymbol{\mu}}_{\kappa}^{\top} \end{aligned}$$

### Beweis.

[Diskrete Verteilung  $(p_1, \dots, p_K)$  der Musterklassen]

Die ML-Zielfunktion lautet zunächst

$$\ell'_{\mathbf{p}}(\omega) = \log \prod_{\kappa=1}^K p_{\kappa}^{T_{\kappa}} = \sum_{\kappa=1}^K T_{\kappa} \log p_{\kappa}$$

und ist aber unter Berücksichtigung der Normierungsbedingung  $\sum_{\kappa} p_{\kappa} = 1$  zu maximieren; die Bedingung wird mit einem Lagrange-Multiplikator inkorporiert:

$$\ell_{\mathbf{p}}(\omega) = \sum_{\kappa=1}^K T_{\kappa} \log p_{\kappa} - \lambda \cdot \left( \sum_{\kappa} p_{\kappa} - 1 \right)$$

Wir bilden nun die partiellen Ableitungen

$$\frac{\partial \ell_{\mathbf{p}}(\omega)}{\partial p_{\kappa}} = T_{\kappa} \frac{1}{p_{\kappa}} - \lambda \quad \text{und} \quad \frac{\partial \ell_{\mathbf{p}}(\omega)}{\partial \lambda} = 1 - \sum_{\kappa} p_{\kappa}$$

Nullsetzen der Ableitungen ergibt

$$\frac{T_{\kappa}}{p_{\kappa}} = \lambda \Rightarrow p_{\kappa} = \frac{T_{\kappa}}{\lambda}$$

und wegen

$$1 = \sum_{\kappa} p_{\kappa} = \sum_{\kappa} \frac{T_{\kappa}}{\lambda} = \frac{1}{\lambda} \sum_{\kappa} T_{\kappa} = \frac{1}{\lambda} \cdot T$$

folgt  $\lambda = T$  und daher  $p_{\kappa} = T_{\kappa}/T$  für alle  $\kappa = 1, \dots, K$ .

## Beweis.

[Parameter  $\mu$  einer univariaten Gaußdichte]

$$f_{\mathbb{X}}(x) = \mathcal{N}(x \mid \mu, \sigma^2) = \frac{1}{\sigma \sqrt{2\pi}} \cdot \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

Die ML-Zielfunktion  $\ell_{\mu, \sigma^2}(\omega) = -2 \cdot \log \prod_{x \in \omega} \mathcal{N}(x \mid \mu, \sigma^2)$  lautet

$$\ell_{\mu, \sigma^2}(\omega) = -2 \cdot \sum_{x \in \omega} \left( -\frac{1}{2} \log(2\pi\sigma^2) - \frac{1}{2} \frac{(x - \mu)^2}{\sigma^2} \right) = T \cdot \log(2\pi\sigma^2) + \frac{1}{\sigma^2} \sum_{x \in \omega} (x - \mu)^2$$

Partielle Ableitung nach  $\mu$ :

$$\frac{\partial \ell(\omega)}{\partial \mu} = \frac{1}{\sigma^2} \sum_{x \in \omega} 2 \cdot (x - \mu) \cdot (-1) = -\frac{2}{\sigma^2} \left( \sum_{x \in \omega} x - \sum_{x \in \omega} \mu \right)$$

Nullsetzen ergibt

$$\sum_{x \in \omega} x = \sum_{x \in \omega} \mu = T \cdot \mu \quad \Rightarrow \quad \hat{\mu} = \frac{1}{T} \sum_{x \in \omega} x$$

□

## Beweis.

[Parameter  $\mu$  einer multivariaten Gaußdichte]

$$\mathcal{N}(x \mid \mu, S) = |2\pi S|^{-1/2} \cdot \exp\left(-\frac{1}{2}(x - \mu)^\top S^{-1}(x - \mu)\right)$$

Die ML-Zielfunktion lautet

$$\begin{aligned} \ell_{\mu, S}(\omega) &= -2 \cdot \log \prod_{x \in \omega} \mathcal{N}(x \mid \mu, S) = -2 \sum_{x \in \omega} \left( -\frac{1}{2} \log |2\pi S| - \frac{1}{2} (x - \mu)^\top S^{-1}(x - \mu) \right) \\ &= T \log |2\pi S| + \sum_{x \in \omega} (x - \mu)^\top S^{-1}(x - \mu) \\ &= T \log |2\pi S| + \sum_{x \in \omega} \left( x^\top S^{-1} x - 2x^\top S^{-1} \mu + \mu^\top S^{-1} \mu \right) \end{aligned}$$

Partielle Ableitung nach  $\mu$  (Gradientenvektor):

$$\begin{aligned} \nabla_{\mu} \ell_{\mu, S}(\omega) &= 0 - 0 + \sum_{x \in \omega} \nabla_{\mu} \left( x^\top S^{-1} x - 2x^\top S^{-1} \mu + \mu^\top S^{-1} \mu \right) \\ &= \sum_{x \in \omega} \left( 0 - 2 \cdot S^{-1} x + 2 \cdot S^{-1} \mu \right) = 2 \cdot S^{-1} \sum_{x \in \omega} (\mu - x) = 2 \cdot S^{-1} \left( T\mu - \sum_{x \in \omega} x \right) \end{aligned}$$

Nullsetzen und Multiplikation mit  $\frac{1}{2} \cdot S$  ergibt

$$T\mu = \sum_{x \in \omega} x \quad \Rightarrow \quad \mu = \frac{1}{T} \sum_{x \in \omega} x$$

□

## Beweis.

[Parameter  $\sigma^2$  einer univariaten Gaußdichte bei bekanntem Wert  $\mu$ ]

$$f_{\mathbb{X}}(x) = \mathcal{N}(x \mid \mu, \sigma^2) = \frac{1}{\sigma \sqrt{2\pi}} \cdot \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

Die ML-Zielfunktion  $\ell_{\mu, \sigma^2}(\omega) = -2 \cdot \log \prod_{x \in \omega} \mathcal{N}(x \mid \mu, \sigma^2)$  lautet

$$\ell_{\mu, \sigma^2}(\omega) = -2 \cdot \sum_{x \in \omega} \left( -\frac{1}{2} \log(2\pi\sigma^2) - \frac{1}{2} \frac{(x - \mu)^2}{\sigma^2} \right) = T \cdot \log(2\pi\sigma^2) + \frac{1}{\sigma^2} \sum_{x \in \omega} (x - \mu)^2$$

Partielle Ableitung nach  $\sigma^2$ :

$$\frac{\partial \ell(\omega)}{\partial \sigma^2} = T \cdot \frac{1}{2\pi\sigma^2} \cdot 2\pi - \frac{1}{\sigma^4} \sum_{x \in \omega} (x - \mu)^2 = \frac{1}{\sigma^2} \left( T - \frac{1}{\sigma^2} \sum_{x \in \omega} (x - \mu)^2 \right)$$

Nullsetzen ergibt

$$T = \frac{1}{\sigma^2} \sum_{x \in \omega} (x - \mu)^2 \quad \Rightarrow \quad \hat{\sigma}^2 = \frac{1}{T} \sum_{x \in \omega} (x - \mu)^2$$

□

### Bemerkung

In der Praxis ist mit  $\sigma^2$  natürlich auch  $\mu$  unbekannt und es muß unter Zuhilfenahme des ML-Schätzwertes  $\hat{\mu}$  optimiert werden. Eine Rechnung ähnlich der obigen ergibt die Varianzschätzformel

$$\hat{\sigma}^2 = \frac{1}{T-1} \sum_{x \in \omega} (x - \hat{\mu})^2.$$

## Beweis.

[Parameter  $S$  einer multivariaten Gaußdichte]

Die ML-Zielfunktion lautet

$$\begin{aligned} \ell_{\mu, S}(\omega) &= T \log |2\pi S| + \sum_{x \in \omega} (x - \mu)^\top S^{-1}(x - \mu) \\ &= TD \log(2\pi) - T \log |S^{-1}| + \sum_{x \in \omega} \text{spur} \left( S^{-1}(x - \mu)(x - \mu)^\top \right) \\ &= TD \log(2\pi) - T \log |S^{-1}| + \underbrace{\text{spur} \left( S^{-1} \cdot \sum_{x \in \omega} (x - \mu)(x - \mu)^\top \right)}_{T \cdot \text{spur}(S^{-1} \cdot \hat{S})} \end{aligned}$$

Wir reformulieren die Zielgröße unter Verwendung der *inversen* Kovarianzmatrix  $Q = S^{-1}$ :

$$\ell_{\mu, Q}(\omega) = TD \log(2\pi) - T \log |Q| + T \cdot \text{spur}(Q \cdot \hat{S})$$

Und nun leiten wir partiell nach der *inversen* Kovarianzmatrix ab:

$$\nabla_Q \ell_{\mu, Q}(\omega) = 0 - T \cdot Q^{-1} + T \cdot \hat{S} = T \cdot (\hat{S} - Q^{-1}) = T \cdot (\hat{S} - S)$$

Nach dem Nullsetzen ergibt sich folglich

$$S = \hat{S} = \frac{1}{T} \sum_{x \in \omega} (x - \mu)(x - \mu)^\top$$

□

## ML-Schätzung für den NV-Klassifikator

Diagonale Kovarianzmatrizen &amp; Mahalanobis-Klassifikator

## Diagonale Kovarianzen

Die ML-Zielgröße zerfällt auf Grund der Unabhängigkeitsannahme in  $(1 + K \cdot D)$  unabhängige Optimierungsterme.

$$\hat{\sigma}_{\kappa,d}^2 = \frac{1}{T_\kappa} \sum_{\mathbf{x} \in \omega_\kappa} (x_d - \mu_{\kappa,d})^2$$

## Mahalanobis-Klassifikator

Bei klassenübergreifenden Kovarianzstatistiken zerfällt  $\ell_\theta(\cdot)$  nicht mehr vollständig in klassenspezifische Optimierungsausdrücke!

$$\hat{\mathbf{S}}_0 = \mathbf{S}_W([\omega_\kappa]) = \frac{1}{T} \sum_{\kappa=1}^K \sum_{\mathbf{x} \in \omega_\kappa} (\mathbf{x} - \boldsymbol{\mu}_\kappa)(\mathbf{x} - \boldsymbol{\mu}_\kappa)^\top$$

Einphasige Berechnung von  $\hat{\mathbf{S}}_0$  ist möglich:  $\mathbf{S}_W = \mathbf{S} - \mathbf{S}_B$

## ML-Schätzung und Lernstichprobenumfang

## Problem

In der NVK-Prüfgröße treten die **Inversen** und die **reziproken Determinanten** aller  $\hat{\mathbf{S}}_\kappa$  auf!

1. Der Varianz-MLS  $\hat{\sigma}_{\kappa,d}$  wird Null, sobald  $|\omega_\kappa| \leq 1$  ist.
2. Der Kovarianz-MLS  $\hat{\mathbf{S}}_\kappa$  wird singulär, sobald  $|\omega_\kappa| \leq D$  ist.
3. Selbst für Klassen mit  $|\omega_\kappa| > D$  besitzt  $\hat{\mathbf{S}}_\kappa$  häufig schlechte Kondition.

➡ **Schwierigkeiten für kleine  $T$ , große  $D$ , große  $K$ .**

## Lösung

Verringerung der **Modellkapazität** (Anzahl freier Parameter)

1. Fixierung und/oder Verklebung von Parametern
2. Strukturierung von Variablenabhängigkeiten
3. Wissensbasierte Engführung des Parameterraums

## ML-Schätzung für den NV-Klassifikator

Richter-Modell: ähnliche Klassenkovarianzen  $\mathbf{S}_\kappa = \alpha_\kappa \mathbf{S}_0$ 

## Iterationsanfang

Berechne Probenstatistiken und initiale Skalierungsfaktoren:

$$\begin{aligned} \hat{p}_\kappa &= \frac{T_\kappa}{T} & \hat{\boldsymbol{\mu}}_\kappa &= \frac{1}{T_\kappa} \sum_{\mathbf{x} \in \omega_\kappa} \mathbf{x} \\ \alpha_\kappa^{(0)} &= 1 & \hat{\mathbf{S}}_\kappa &= \frac{1}{T_\kappa} \sum_{\mathbf{x} \in \omega_\kappa} \mathbf{x} \mathbf{x}^\top - \hat{\boldsymbol{\mu}}_\kappa \hat{\boldsymbol{\mu}}_\kappa^\top \end{aligned}$$

## Iterationsschritt

Berechne Kovarianzprototyp und Skalierungsfaktoren für  $i = 1, 2, \dots$ :

$$\begin{aligned} \mathbf{s}_0^{(i)} &= \sum_{\kappa=1}^K \hat{p}_\kappa \cdot (\alpha_\kappa^{(i-1)})^{-1} \cdot \hat{\mathbf{S}}_\kappa \\ \alpha_\kappa^{(i)} &= \frac{1}{D} \cdot \text{spur} \left( \hat{\mathbf{S}}_\kappa \cdot (\mathbf{s}_0^{(i)})^{-1} \right) \end{aligned}$$

Multivariate Normalverteilungsdichte

Normalverteilungsklassifikatoren

Maximum-Likelihood Parameterschätzung

Maximum-a posteriori- und Bayesschätzung

Graphische Gaußsche Modelle

Mathematische Hilfsmittel



## Maximum-a posteriori Schätzung

Verteilungsparameter  $\theta$  als Werte einer Zufallsvariablen  $\Theta$

### Bayesscher Denkansatz

Die wahren Verteilungsparameter  $\theta^*$  des Prozesses sind nicht nur **unbekannt**, sie sind sogar **stochastisch**.

Ihre Verteilungsdichte  $f_{\Theta}(\cdot)$  repräsentiert unser **Vorwissen** über ihre möglichen Werte(kombinationen).

### Lemma

Sind die Parameter der Verteilungsfamilie  $\{f_{\mathbb{X}}(\cdot|\theta)\}_{\theta \in \mathcal{M}}$  selbst gemäß **a priori Dichte**  $f_{\Theta}(\theta)$  verteilt, so lautet — für den unabhängigen und identisch gezogenen Datensatz  $\omega$  — die datenbedingte **a posteriori Dichte** der Parameter

$$P(\theta|\omega) = \frac{P(\theta) \cdot P(\omega|\theta)}{P(\omega)} = \frac{f_{\Theta}(\theta) \cdot \prod_{\mathbf{x} \in \omega} f_{\mathbb{X}}(\mathbf{x}|\theta)}{P(\omega)}.$$

## Wissenswertes über die Maximum-a posteriori Schätzung

### Spezialfall Maximum-Likelihood

Unter Gleichverteilungsannahme für  $f_{\Theta}(\cdot)$  mutiert die MAP-Schätzung in eine ML-Schätzung.

### Asymptotisches Schätzverhalten

Für große Stichproben ( $|\omega| \rightarrow \infty$ ) strebt  $\hat{\theta}_{\text{MAP}}$  gegen  $\hat{\theta}_{\text{ML}}$ .

### Methode der konjugierten Dichtefamilien

Die *analytische* Optimierung der MAP-Zielfunktion erfordert eine geeignete Form der a priori-Dichte:

$$f_{\Theta}(\theta) \hat{=} C \cdot \prod_{\mathbf{z} \in \omega_{\text{prior}}} f_{\mathbb{X}}(\mathbf{z}|\theta)$$

Mit dieser Wahl gilt nämlich

$$\hat{\theta}_{\text{MAP}}(\omega) = \hat{\theta}_{\text{ML}}(\omega \cup \omega_{\text{prior}})$$

und das Problem der  $f_{\Theta}(\cdot)$ -Findung ist auf elegante Art gelöst!

## Maximum-a posteriori Schätzung

Die im Lichte der Datenprobe wahrscheinlichsten Verteilungsparameter

### Definition

Die **Maximum-a posteriori-Schätzung** (MAP) der Parameter einer Dichtefamilie  $[f(\mathbf{x}|\theta)]$  unter Annahme der **a priori**-Verteilungsdichte  $f_{\Theta}(\theta)$  für  $\theta$  maximiert die stichprobenbedingte Wahrscheinlichkeit des gesuchten Parameterfeldes, d.h. es gilt:

$$\hat{\theta}_{\text{MAP}} = \underset{\theta}{\operatorname{argmax}} \left( f_{\Theta}(\theta) \cdot \prod_{\mathbf{x} \in \omega} f_{\mathbb{X}}(\mathbf{x}|\theta) \right)$$

### Bemerkungen

1. Der MAP-Schätzwert  $\hat{\theta}_{\text{MAP}}$  ist von allen Parameterwerten derjenige, der zu den vorliegenden Daten  $\omega$  am besten paßt.
2. Hand aufs Herz — niemand (außer dem *Capo di tutti capi*) kennt diese mysteriöse Dichte  $f_{\Theta}(\cdot)$ .

## MAP-Schätzung für diskrete Verteilungen

Wahrscheinlichkeitsparameter  $p_1 + p_2 + \dots + p_K = 1$  für  $K$  Ereignisse

### Definition

Der Zufallsvektor  $\Theta = (\Theta_1, \dots, \Theta_K)^T \in [0, 1]^K$  mit  $\sum_{\ell} \Theta_{\ell} = 1$  heißt **Dirichlet-verteilt** mit den **Hyperparametern**  $r_1, \dots, r_K > -1$  genau dann, wenn gilt:

$$f_{\Theta}(\mathbf{p}) = \mathcal{D}(\mathbf{p}|\mathbf{r}) = C \cdot \prod_{\ell=1}^K p_{\ell}^{r_{\ell}}$$

### Bemerkungen

1. Für  $\mathbf{r} = \mathbf{0}$  ist  $\mathcal{D}(\mathbf{p}|\mathbf{r})$  eine **Gleichverteilung**.
2. Für  $\mathbf{r} = \mathbf{1}$  nimmt  $\mathcal{D}(\mathbf{p}|\mathbf{r})$  ihr Dichtemaximum bei der **Gleichverteilung**  $p_{\ell} \equiv 1/K$  an.
3. Allgemein nimmt  $\mathcal{D}(\mathbf{p}|\mathbf{r})$  ihr Dichtemaximum bei der Verteilung  $\mathbf{p} \propto \mathbf{r}$  an, also für die Wahrscheinlichkeiten  $p_{\ell} = r_{\ell}/R$ ,  $R = \sum_i r_i$ .
4. Der Dichtegipfel ist umso steiler, je größer der Skalenfaktor  $R$  ist.

## MAP-Schätzung für diskrete Verteilungen

### Satz

Gehorchen die kanonischen Parameter  $p_1, \dots, p_K$  einer diskreten Wahrscheinlichkeitsverteilung der Dirichletverteilung mit Hyperparametern  $\mathbf{r} \in \mathbb{R}^K$ , so lautet der MAP-Schätzwert für eine Stichprobe mit den absoluten Ereignishäufigkeiten  $T_1 + T_2 + \dots + T_K = T$

$$\hat{p}_\ell = \frac{T_\ell + r_\ell}{T + R}, \quad R = \sum_{\ell=1}^K r_\ell.$$

### Bemerkungen

1. Eine MAP-Schätzung mit Vorwissen  $\mathcal{D}(\cdot|\mathbf{r})$  bewirkt die Aufstockung der Lerndaten  $\omega$  um eine **virtuelle Datenprobe**  $\omega_{\text{prior}}$  mit den Ereignishäufigkeiten  $r_\ell$ ; diese Werte müssen allerdings nicht unbedingt ganzzahlig sein.
2. Der Spezialfall einer gleichverteilten oder **uninformativen** Dirichletdichte ( $r_\ell \equiv r_0$ ) ergibt die MAP-Schätzwerte (**Laplaceschätzformel** im Fall  $r_0 = 1$ )

$$\hat{p}_\ell = \frac{(T_\ell + r_0)}{(T + K \cdot r_0)}, \quad \ell = 1, 2, \dots, K.$$

## MAP-Schätzung für die multivariate NV-Dichte

### Definition

Eine Zufallsmatrix  $\mathbf{S}$  über der Mannigfaltigkeit aller symmetrischen, positiv-definiten  $(D \times D)$ -Matrizen heißt **Wishart-verteilt** genau dann, wenn

$$f_{\mathbf{S}}(\mathbf{S}) = \mathcal{W}(\mathbf{S} | \alpha, \mathbf{V}) = \frac{1}{2^{\frac{\alpha D}{2}} |\mathbf{V}|^{\frac{\alpha}{2}} \Gamma_D(\frac{\alpha}{2})} \cdot |\mathbf{S}|^{\frac{\alpha-D-1}{2}} \cdot \exp(-\frac{1}{2} \cdot \text{spur}(\mathbf{V}^{-1} \mathbf{S}))$$

gilt mit den Hyperparametern  $\alpha > D - 1$  und  $\mathbf{V} \in \mathbb{R}^{D \times D}$  positiv-definit.

### Lemma

Für die multivariate NV-Dichte  $\mathcal{N}(\boldsymbol{\mu}, \mathbf{S})$  bildet das Produkt

$$f_{\Theta}(\boldsymbol{\mu}, \mathbf{S}) = \mathcal{N}(\boldsymbol{\mu} | \mathbf{m}, \tau^{-1} \mathbf{S}) \cdot \mathcal{W}(\mathbf{S}^{-1} | \alpha, \mathbf{V})$$

eine konjugierte Dichtefamilie mit den Hyperparametern  $\mathbf{m} \in \mathbb{R}^D$ ,  $\tau > 0$ ,  $\alpha > D - 1$  und positiv-definiten Matrix  $\mathbf{V} \in \mathbb{R}^{D \times D}$ .

### Beweis.

Es beträgt die Stichprobenwahrscheinlichkeit

$$P(\omega | \mathbf{p}) = \prod_{\kappa=1}^K p_{\kappa}^{T_{\kappa}}$$

und die a posteriori Parameterwahrscheinlichkeit (bei festen Hyperparametern)

$$P(\mathbf{p} | \omega) \propto P(\omega | \mathbf{p}) \cdot f_{\Theta}(\mathbf{p}) \propto \prod_{\kappa=1}^K p_{\kappa}^{T_{\kappa}} \cdot \prod_{\kappa=1}^K p_{\kappa}^{r_{\kappa}} \propto \prod_{\kappa=1}^K p_{\kappa}^{(T_{\kappa} + r_{\kappa})}$$

Das Maximum nimmt  $P(\mathbf{p} | \omega)$  bekanntlich für diejenige Verteilung an, die proportional zu den Exponenten ist:

$$\hat{p}_{\kappa} = \frac{T_{\kappa} + r_{\kappa}}{T + R}, \quad R = \sum_{\kappa} r_{\kappa}$$

□

Der MAP-Schätzwert ist ein gewichtetes Mittel („Konvexkombination“) aus ML-Schätzwert und dem Modus

$$\rho_{\kappa} = r_{\kappa} / R, \quad \kappa = 1, \dots, K$$

der a priori-Dichte:

$$\hat{p}_{\kappa} = \frac{T_{\kappa} + r_{\kappa}}{T + R} = \frac{T_{\kappa}}{T + R} + \frac{r_{\kappa}}{T + R} = \underbrace{\frac{T_{\kappa}}{T}}_{\hat{p}_{\kappa}^{\text{ML}}} \cdot \underbrace{\frac{T}{T + R}}_{\lambda} + \underbrace{\frac{r_{\kappa}}{R}}_{\rho_{\kappa}} \cdot \underbrace{\frac{R}{T + R}}_{(1-\lambda)}$$

## MAP-Schätzung für den NV-Klassifikator

### Satz

Die Lerndaten  $\omega_1, \dots, \omega_K \subset \mathbb{R}^D$  eines numerischen Klassifikationsproblems seien klassenweise normalverteilt mit den unbekannten Parametern  $(p_{\kappa}, \boldsymbol{\mu}_{\kappa}, \mathbf{S}_{\kappa})$ ,  $\kappa = 1, \dots, K$ . Die a priori Verteilung der Parameter sei definiert durch

$$f_{\Theta}(\boldsymbol{\theta}) = \mathcal{D}(\mathbf{p} | \mathbf{r}) \cdot \prod_{\kappa=1}^K \mathcal{N}(\boldsymbol{\mu}_{\kappa} | \mathbf{m}_{\kappa}, \tau_{\kappa}^{-1} \mathbf{S}_{\kappa}) \cdot \prod_{\kappa=1}^K \mathcal{W}(\mathbf{S}_{\kappa}^{-1} | \alpha_{\kappa}, \mathbf{V}_{\kappa}).$$

Dann lauten die Maximum-a posteriori-Parameter:

$$\begin{aligned} \hat{p}_{\kappa} &= \frac{r_{\kappa} + T_{\kappa}}{R + T}, \quad R = \sum_{\kappa} r_{\kappa} \\ \hat{\boldsymbol{\mu}}_{\kappa} &= \frac{1}{T_{\kappa} + T_{\kappa}} \left( \tau_{\kappa} \mathbf{m}_{\kappa} + \sum_{\mathbf{x} \in \omega_{\kappa}} \mathbf{x} \right) \\ \hat{\mathbf{S}}_{\kappa} &= \frac{\mathbf{V}_{\kappa} + \tau_{\kappa} (\hat{\boldsymbol{\mu}}_{\kappa} - \mathbf{m}_{\kappa})(\hat{\boldsymbol{\mu}}_{\kappa} - \mathbf{m}_{\kappa})^{\top} + \sum_{\mathbf{x} \in \omega_{\kappa}} \mathbf{x} \mathbf{x}^{\top} - T_{\kappa} \hat{\boldsymbol{\mu}}_{\kappa} \hat{\boldsymbol{\mu}}_{\kappa}^{\top}}{(\alpha_{\kappa} - D) + T_{\kappa}} \end{aligned}$$

## „Plug-in“-Schätzverfahren

Die Suche nach den *unbekannten*, aber *wahren* Parametern

### Traditionelles Induktionsparadigma

Die Verteilungsannahme  $\omega \sim f_{\mathbb{X}}(\cdot|\theta)$  ist korrekt.

Es existiert eine wahre Parameterkonfiguration  $\theta^*$  — wir müssen sie nur finden!

### ML-Schätzung

$$\hat{\theta}_{\text{ML}} = \underset{\theta}{\operatorname{argmax}} P(\omega|\theta)$$

### Posterior-Mean-Schätzung

$$\hat{\theta}_{\text{PM}} = \mathcal{E}[\Theta|\omega] = \int \theta \cdot P(\theta|\omega) d\theta$$

### MAP-Schätzung

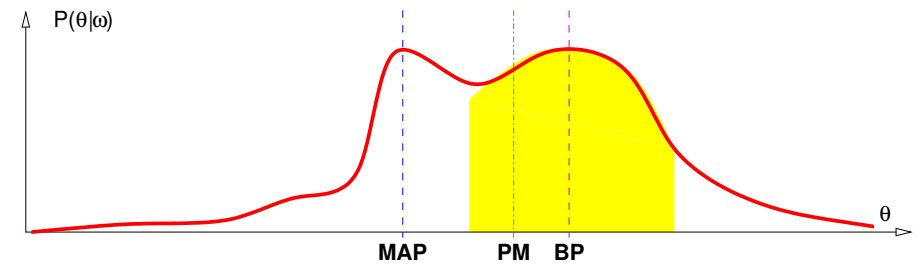
$$\hat{\theta}_{\text{MAP}} = \underset{\theta}{\operatorname{argmax}} P(\theta|\omega)$$

### Bayespunkt-Schätzung

$$\hat{\theta}_{\text{BP}}^{(\rho)} = \underset{\theta}{\operatorname{argmax}} \int_{\mathcal{U}_{\rho}(\theta)} P(\vartheta|\omega) d\vartheta$$

## „Plug-in“-Schätzverfahren

Analyse der a posteriori Parameterdichte



**MAP** Wo liegt der **Gipfel** der Posterioridichte?

**PM** Wo liegt der **Durchschnitt** der Posterioridichte?

**BP** Wo liegt das **kleinste Intervall** mit Wahrscheinlichkeitsmasse  $\rho > 0$ ?

## Bayes-Schätzung

Der Abschied von der Idee „wahrer“ Verteilungsparameter

### Bayessches Induktionsparadigma

Die Verteilungsannahme  $\omega \sim f_{\mathbb{X}}(\cdot|\theta)$  ist korrekt.

Aber jedes  $\mathbf{x} \in \omega$  wird unter Verwendung eines eigenen, zufällig ausgewürfelten Modellparameters  $\theta$  gezogen!

$$\begin{aligned} P(\mathbf{x}|\omega) &= \int_{\mathcal{M}} P(\mathbf{x}, \boldsymbol{\theta} | \omega) d\boldsymbol{\theta} \\ &= \int_{\mathcal{M}} P(\mathbf{x} | \boldsymbol{\theta}, \omega) \cdot P(\boldsymbol{\theta} | \omega) d\boldsymbol{\theta} \\ &= \int_{\mathcal{M}} \underbrace{f_{\mathbb{X}}(\mathbf{x}|\boldsymbol{\theta})}_{\text{Modellidichte}} \cdot \underbrace{\frac{f_{\mathbb{X}}(\omega|\boldsymbol{\theta}) \cdot f_{\Theta}(\boldsymbol{\theta})}{f_{\mathbb{X}}(\omega)}}_{\text{a posteriori}} d\boldsymbol{\theta} \end{aligned}$$

Analytisch extrem schwer lösbar — bestenfalls wenn  $f_{\Theta}(\cdot) \equiv c$

## Bayesapproximation

Asymptotisch korrekte Näherung unter Gleichverteilungsannahme für  $f_{\Theta}(\cdot)$

### Praktikable Näherungslösung für den Bayesschätzer

Unwissen um  $f_{\Theta}(\cdot) \rightsquigarrow$  Gleichverteilung  $\rightsquigarrow$  Herauskürzen

Simultan in Zähler und Nenner: Integralbildung  $\rightsquigarrow$  Maximumbildung

$$\begin{aligned} P(\mathbf{x}|\omega) &= \frac{P(\mathbf{x}, \omega)}{P(\omega)} = \frac{\int f_{\mathbb{X}}(\omega, \mathbf{x}|\boldsymbol{\theta}) \cdot f_{\Theta}(\boldsymbol{\theta}) d\boldsymbol{\theta}}{\int f_{\mathbb{X}}(\omega|\boldsymbol{\theta}) \cdot f_{\Theta}(\boldsymbol{\theta}) d\boldsymbol{\theta}} \\ &\approx \frac{\max_{\boldsymbol{\theta}} f_{\mathbb{X}}(\omega, \mathbf{x}|\boldsymbol{\theta})}{\max_{\boldsymbol{\theta}} f_{\mathbb{X}}(\omega|\boldsymbol{\theta})} = \frac{\prod_{\mathbf{z} \in \omega, \mathbf{x}} f_{\mathbb{X}}(\mathbf{z} | \hat{\boldsymbol{\theta}}_{\text{ML}}(\omega, \mathbf{x}))}{\prod_{\mathbf{z} \in \omega} f_{\mathbb{X}}(\mathbf{z} | \hat{\boldsymbol{\theta}}_{\text{ML}}(\omega))} \end{aligned}$$

**Achtung:**

Die Bayesapproximation  $\hat{P}_{\text{BA}}(\mathbf{x}|\omega)$  ist i.a. **keine** Dichtefunktion (Normierungseigenschaft)!

Multivariate Normalverteilungsdichte

Normalverteilungsklassifikatoren

Maximum-Likelihood Parameterschätzung

Maximum-a posteriori- und Bayesschätzung

Graphische Gaußsche Modelle

Mathematische Hilfsmittel

## Graphische Gaußsche Modelle

Die Bias-Varianz-Problematik

## Dichtemodell mit vielen Parametern

NV-Dichten mit voll besetzter Kovarianzmatrix

Alle paarweisen Merkmalabhängigkeiten  $\rightsquigarrow O(KD^2)$ 

Kleiner Bias — große Varianz

## Dichtemodell mit wenigen Parametern

NV-Dichten mit diagonal besetzter Kovarianzmatrix

Alle Merkmale paarweise unabhängig  $\rightsquigarrow O(KD)$ 

Großer Bias — kleine Varianz

## Lösungsidee

Nicht alle, sondern nur die **wichtigen** Merkmalabhängigkeiten werden explizit modelliert.

## Gaußsche Bayesnetze

## Kettenregel der Wahrscheinlichkeitstheorie

$$P(x_1, \dots, x_D) = P(x_1) \cdot P(x_2 | x_1) \cdot \prod_{d=3}^D P(x_d | x_1, \dots, x_{d-1})$$

Das  $d$ -te Merkmal ist explizit von  $(d-1)$  anderen abhängig.

Beispiel: baumförmige Bayesnetze

$$P(x_1, \dots, x_D) \approx \prod_{d=1}^D P(x_d | x_{\pi(d)})$$

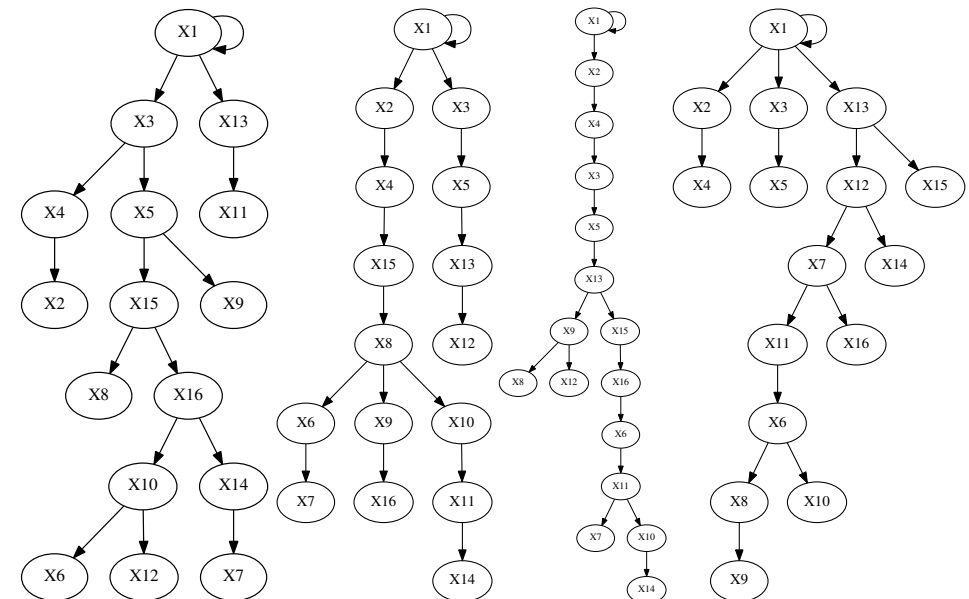
Jedes Merkmal  $x_d$  ist explizit nur von **genau einem** anderen abhängig.

## Problem

Finde diejenige Abhängigkeitsstruktur, welche die exakteste Näherung der Datenverteilung gewährleistet!

## Gaußsche Bayesnetze

Datensatz letter.lern (16 Merkmale, Klassen 'A', 'B', 'C', 'D')



## Gaußsche Markovnetze

### Parametrische Struktur der multivariaten NV-Dichte

$$-2 \cdot \log \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}, \mathbf{S}) = |2\pi\mathbf{S}| + \sum_{i=1}^D \sum_{j=1}^D (x_i - \mu_i) \cdot C_{ij} \cdot (x_j - \mu_j), \quad \mathbf{C} := \mathbf{S}^{-1}$$

Modellkomplexität  $\hat{=}$  Anzahl nicht verschwindender Einträge von  $\mathbf{S}^{-1}$

### Aufgabenstellung der Kovarianzselektion

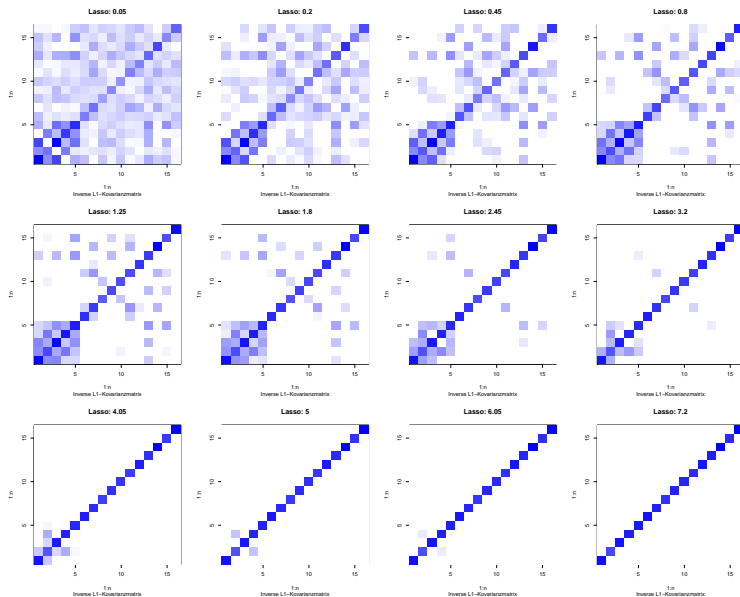
Suche eine Näherungsmatrix  $\tilde{\mathbf{S}} \approx \hat{\mathbf{S}}$ , deren Inverse möglich **viele Nulleinträge** aufweist!

### Bedingte statistische Unabhängigkeit

Über normalverteilte Daten wissen wir, daß  $C_{ij} = 0$  genau dann gilt, wenn die beiden Merkmale  $x_i$  und  $x_j$  **statistisch unabhängig** sind, sofern wir die Kenntnis der restlichen Merkmale  $\{x_1, \dots, x_D\} \setminus \{x_i, x_j\}$  voraussetzen.

## Gaußsche Markovnetze

### Lasso (regularisierte $\|\cdot\|_1$ -Norm Matrixinvertierung)



### Beispiel

Datensatz  
letter  
16 Merkmale  
alle Klassen

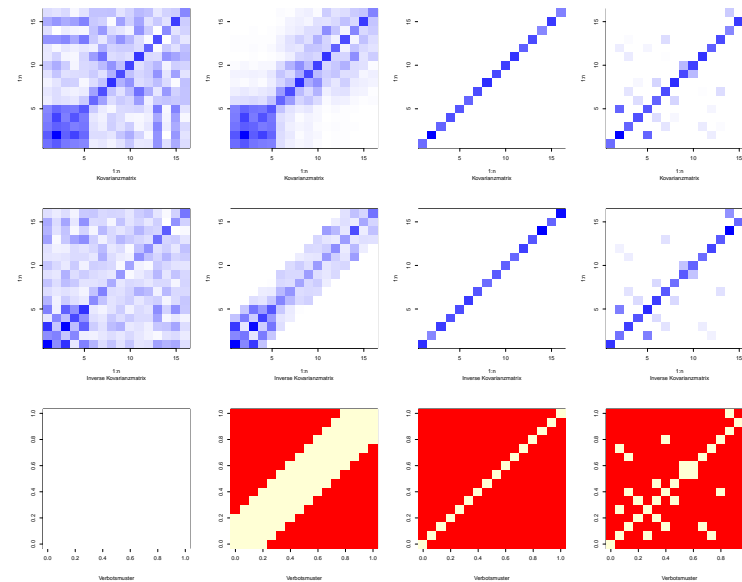
Konzentrations-  
matrizen für  
unterschiedliche  
**Regularisierungs-  
parameter**

$$\rho = \frac{1}{20} \cdot n^2$$

$$n = 1, 2, \dots, 12$$

## Gaußsche Markovnetze

### Dempsters Kovarianzselektion



### Beispiel

Datensatz  
letter  
16 Merkmale  
alle Klassen

*oben:*  
Kovarianz  
 $\hat{\mathbf{S}} = \mathbf{C}^{-1}$

*Mitte:*  
Konzentration  
 $\mathbf{C}$  erfüllt  $\mathbf{A}$

*unten:*  
Adjazenz  $\mathbf{A}$   
Abhängigkeits-  
muster  
(gegeben)

Multivariate Normalverteilungsdichte

Normalverteilungsklassifikatoren

Maximum-Likelihood Parameterschätzung

Maximum-a posteriori- und Bayesschätzung

Graphische Gaußsche Modelle

Mathematische Hilfsmittel

## Gradientenvektor und Gradientenmatrix

Extremalwertaufgabe  $\Rightarrow$  „Ableiten & Nullsetzen“

### Definition

Es sei

$$f: \mathbb{R}^n \rightarrow \mathbb{R} \quad \text{bzw.} \quad g: \mathbb{R}^{n \times m} \rightarrow \mathbb{R}$$

ein Vektor- bzw. ein Matrixfunktional. Dann heißen die Felder

$$\nabla_{\mathbf{x}} f(\mathbf{x}) = \begin{pmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{pmatrix} \quad \text{bzw.} \quad \nabla_{\mathbf{Y}} f(\mathbf{Y}) = \begin{pmatrix} \frac{\partial g}{\partial y_{11}} & \frac{\partial g}{\partial y_{12}} & \cdots & \frac{\partial g}{\partial y_{1m}} \\ \frac{\partial g}{\partial y_{21}} & \frac{\partial g}{\partial y_{22}} & \cdots & \frac{\partial g}{\partial y_{2m}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial g}{\partial y_{n1}} & \frac{\partial g}{\partial y_{n2}} & \cdots & \frac{\partial g}{\partial y_{nm}} \end{pmatrix}$$

von partiellen Ableitungen nach allen Eingangsvariablen der

**Gradientenvektor** von  $f$  an der Stelle  $\mathbf{x}$  bzw. die **Gradientenmatrix** von  $g$  an der Stelle  $\mathbf{Y}$ .

### Bemerkung

Eine notwendige Bedingung für das Vorliegen eines relativen Maximums von  $f$  an der Stelle  $\mathbf{x} \in \mathbb{R}^n$  ist das Verschwinden ( $\nabla_{\mathbf{x}} f = \mathbf{0}$ ) des Gradientenvektors in diesem Punkt.

## Gradientenvektorberechnung

Beispiel: Quadratische Form

### Beispielfunktional

$$f(\mathbf{x}) = \sum_{i=1}^n \sum_{j=1}^n x_i A_{ij} x_j = \mathbf{x}^\top \mathbf{A} \mathbf{x}$$

### Berechnung partieller Ableitungen

$$\frac{\partial f(\mathbf{x})}{\partial x_k} = \sum_{i=1}^n \sum_{j=1}^n \frac{\partial}{\partial x_k} (x_i A_{ij} x_j) = \sum_{i \neq k} x_i A_{ik} + \sum_{j \neq k} A_{kj} x_j + 2x_k \cdot A_{kk} = 2 \cdot \sum_{i=1}^n x_i A_{ik} = 2 \cdot \mathbf{a}_k^\top \mathbf{x}$$

### Gradientenvektor

$$\nabla_{\mathbf{x}} (\mathbf{x}^\top \mathbf{A} \mathbf{x}) = 2 \cdot \mathbf{A} \mathbf{x}$$

## Gradientenvektorberechnung

Beispiel: Linearkombination

### Beispielfunktional

$$f(\mathbf{x}) = \sum_{i=1}^n a_i x_i = \mathbf{a}^\top \mathbf{x}$$

### Berechnung partieller Ableitungen

$$\frac{\partial f(\mathbf{x})}{\partial x_k} = \frac{1}{\partial x_k} \left( \sum_{i=1}^n a_i x_i \right) = \sum_{i=1}^n \frac{\partial}{\partial x_k} (a_i x_i) = \sum_{i=1}^n a_i \cdot \frac{\partial}{\partial x_k} (x_i) = a_k$$

### Gradientenvektor

$$\nabla_{\mathbf{x}} (\mathbf{a}^\top \mathbf{x}) = \mathbf{a}$$

## Gradientenmatrixberechnung

Beispiel: Frobeniusnorm

### Beispielfunktional

$$f(\mathbf{X}) = \sum_{i=1}^n \sum_{j=1}^m X_{ij}^2 = \text{spur} (\mathbf{X}^\top \mathbf{X}) = \|\mathbf{X}\|_{\text{Frob}}^2$$

### Berechnung partieller Ableitungen

$$\frac{\partial f(\mathbf{X})}{\partial X_{k\ell}} = \sum_{i=1}^n \sum_{j=1}^m \frac{\partial}{\partial X_{k\ell}} (X_{ij}^2) = 2 \cdot X_{k\ell}$$

### Gradientenmatrix

$$\nabla_{\mathbf{X}} (\|\mathbf{X}\|_{\text{Frob}}^2) = 2 \cdot \mathbf{X}$$

## Gradientenmatrixberechnung

Beispiel: Determinante

### Beispielfunktional

$$f(\mathbf{X}) = \det(\mathbf{X})$$

### Berechnung partieller Ableitungen

$$\frac{\partial f(\mathbf{X})}{\partial X_{kl}} = \frac{\partial \det(\mathbf{X})}{\partial X_{kl}} = \det(\mathbf{X}) \cdot (\mathbf{X}^{-1})_{kl}$$

### Gradientenmatrix

$$\nabla_{\mathbf{X}} (\det(\mathbf{X})) = \det(\mathbf{X}) \cdot \mathbf{X}^{-1}$$

## Zusammenfassung (7)

1. Die **multivariate Normalverteilung** beschreibt eine unimodale (Zentrum  $\mu$ ), exponentiell abklingende Dichte mit elliptisch-symmetrischen (Trägheitsachsen von  $\mathbf{S}$ ) Isolinien.
2. Die Prüfgrößen der NV-Bayesregel sind **quadratische Polynome** in den Merkmalen  $x_1, \dots, x_D$ .
3. Die **Maximum-Likelihood**-Schätzung sucht die Modellparameter mit der größten Datenerzeugungswahrscheinlichkeit.
4. Die ML-Zielgröße ist nach allen Parametern **partiell abzuleiten**; nach **Nullsetzen der Gradienten** ergibt sich günstigenfalls eine **geschlossene Lösung** (LGS) oder wenigstens eine rasch konvergierende **Iterationsformel**.
5. Die **Maximum-a posteriori**-Schätzung verwendet **a priori-Wissen** über die Dichteparameter und ist **robuster** bei (zu) **kleinen Lernenstichproben**.
6. Praktikable MAP-Schätzer bedienen sich der Methode der **konjugierten Parameterdichtefamilien**.
7. Verteilungsmodelle werden robuster, wenn die **Abhängigkeitsstruktur der Merkmale** sachgemäß **ausgedünnt** wird.

## Gradientenmatrixberechnung

Beispiel: Logarithmierte Determinante

### Beispielfunktional

$$f(\mathbf{X}) = \log \det(\mathbf{X})$$

### Berechnung partieller Ableitungen

$$\frac{\partial f(\mathbf{X})}{\partial X_{kl}} = \frac{\partial \log \det(\mathbf{X})}{\partial \det(\mathbf{X})} \cdot \frac{\partial \det(\mathbf{X})}{\partial X_{kl}} = \frac{1}{\det(\mathbf{X})} \cdot \det(\mathbf{X}) \cdot (\mathbf{X}^{-1})_{kl} = (\mathbf{X}^{-1})_{kl}$$

### Gradientenmatrix

$$\nabla_{\mathbf{X}} (\log \det(\mathbf{X})) = \mathbf{X}^{-1}$$