

# Kapitel 9

## Polynomklassifikator

# Ansätze zum Entwurf eines Klassifikators

Es gibt zwei prinzipiell unterschiedliche Ansätze zum Entwurf eines Klassifikators:

1. Statistische parametrische Modellierung der Klassenverteilungen, dann MAP
2. Lösung eines Abbildungsproblems durch Funktionsapproximation (nichtlineare Regression)

$$\min E \{ \|\mathbf{f}(\mathbf{x}) - \mathbf{y}\|^2 \}$$

$\mathcal{X}$  – Merkmalsraum

$\mathcal{Y}$  – Entscheidungsraum

Zu 1.) Die bisher beschriebene Vorgehensweise beim Entwurf eines Klassifikators beruht darauf, dass man die *klassenspezifischen Verteilungsdichten*

$$p(\mathbf{x}|\omega)$$

parametrisch durch statistische Modelle annähert (durch Schätzung der Parameter, z.Bsp. einer Gauß-Verteilung) und durch eine Maximumselektion zu einer Entscheidung kommt. Lernen bedeutet hierbei: Verbesserung des Parameter-Fitting.

Zu 2.) Es gibt nun eine zweite Möglichkeit der Herangehensweise, welche auf die Auswertung der a-posteriori-Wahrscheinlichkeitsdichte

$$p(\omega|\mathbf{x})$$

aufbaut und durch ein Problem der *Funktionsapproximation* beschrieben werden kann.

Diese Funktionsapproximation kann z.Bsp. durch eine *nichtlineare Regression mit Polynomen* oder auch mit Hilfe eines *künstlichen Neuronalen Netzwerkes* (NN) durchgeführt werden. Im Rahmen dieser Veranstaltung sollen die Grundlagen für beide Vorgehensweisen behandelt werden.

Grundsätzlich ist die Suche nach einer besten Näherungsfunktion ein *Variationsproblem*, welches durch die Wahl von Basisfunktionen auf ein *parametrisches Optimierungsproblem* zurückgeführt werden kann. Lernen bedeutet dann auch hier: Parameterfitting.

Die Gleichwertigkeit der beiden genannten Vorgehensweisen ergibt sich aus dem Bayes-Theorem:

$$p(\omega | \mathbf{x}) = \frac{p(\mathbf{x} | \omega) p(\omega)}{p(\mathbf{x})}$$

$\swarrow$  1/K  
 $\nwarrow$  unabhängig von  $\omega$

Die Gleichwertigkeit ergibt sich daraus, dass der Nenner unabhängig von  $\omega$  ist und die a-priori-W. i.a.  $1/K$  gesetzt werden kann.

Im folgenden soll nun die Überführung in ein Funktionsapproximationsproblem begründet werden.

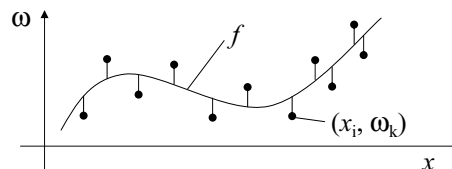
Bei bekannter a-posteriori-W.  $p(\omega|\mathbf{x})$  würde für jedes kontinuierliche  $\mathbf{x}$  bestmöglich ein  $\omega$  zugeordnet werde (funktionale Zuordnung  $f:\mathbf{x} \rightarrow \omega$ ). Gegeben sind jedoch nur *Stichproben* und man sucht nach einer Funktion  $f$ , welche bestmöglich die Einzelexperimente fittet und damit eine Abbildung realisiert:

$$f: \underset{\mathbf{x}}{\text{Merkmalsraum}} \rightarrow \underset{\omega}{\text{Bedeutungsraum}}$$

Diese Aufgabenstellung kann mit Hilfe der Variationsrechnung gelöst werden.

Wählt man als Gütekriterium das minimale Fehlerquadrat, so geht es um die Minimierung von:

$$J = \min_{f(\mathbf{x})} E\{\|\mathbf{f}(\mathbf{x}) - \mathbf{y}\|^2\}$$



Dabei entstehen die Zielvektoren  $\{\mathbf{y}_i\}$  im *Entscheidungsraum*  $\mathcal{Y}$  durch einfache Abbildung der skalaren Bedeutungswerte  $\{\omega_i\}$

$$\Omega := \{\omega_1, \omega_2, \dots, \omega_K\}$$

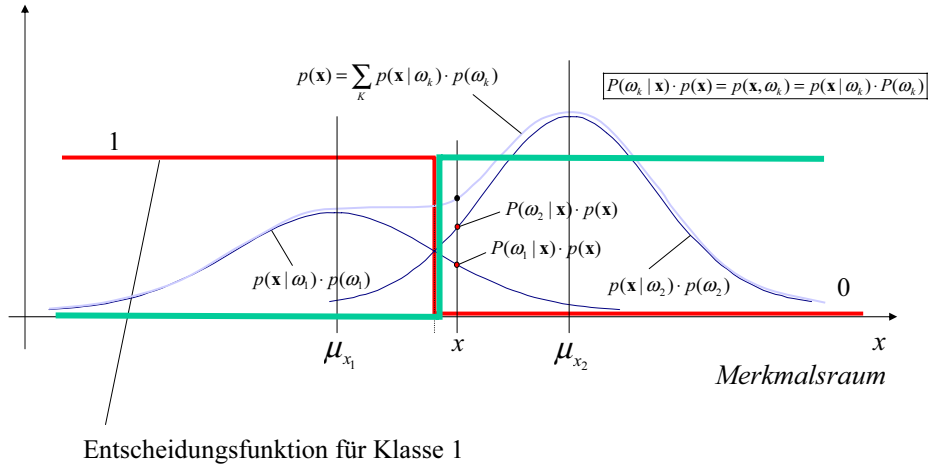


$$\mathcal{Y} := \{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_K\}$$

mit:

$$\mathbf{y}_i = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad \text{i-ter Einheitsvektor}$$

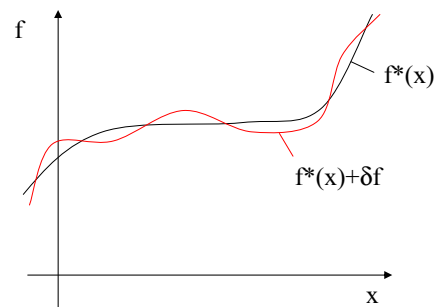
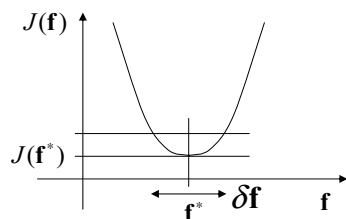
## Zwei-Klassen-Problem mit Gaußverteiligdichte



## Lösung des Variationsproblems

Die Lösung des Variationsproblems erhält man über die Annahme, dass man eine optimale Funktion  $\mathbf{f}^*(\mathbf{x})$  bereits kennt. Dann muss jede Variation zu einer Verschlechterung des Gütekriteriums führen:

$$J(\mathbf{f}^* + \delta\mathbf{f}) > J(\mathbf{f}^*) \quad \text{für } \forall \delta\mathbf{f} \neq \mathbf{0}$$



Mit Hilfe der Definitionsgleichung ergibt sich:

$$\begin{aligned} J(\mathbf{f}^* + \delta\mathbf{f}) &= E \left\{ \|\mathbf{f}^* + \delta\mathbf{f} - \mathbf{y}\|^2 \right\} = E \left\{ \langle (\mathbf{f}^* - \mathbf{y}) + \delta\mathbf{f}, (\mathbf{f}^* - \mathbf{y}) + \delta\mathbf{f} \rangle \right\} \\ &= E \left\{ \|\mathbf{f}^* - \mathbf{y}\|^2 \right\} + 2E \left\{ \delta\mathbf{f}^T (\mathbf{f}^* - \mathbf{y}) \right\} + E \left\{ \|\delta\mathbf{f}\|^2 \right\} \\ &\Rightarrow J(\mathbf{f}^*) = E \left\{ \|\mathbf{f}^* - \mathbf{y}\|^2 \right\} \quad (\text{für } \delta\mathbf{f} = \mathbf{0}) \end{aligned}$$

Eingesetzt in die Ungleichung ergibt:

$$E \left\{ \|\delta\mathbf{f}\|^2 \right\} + 2E \left\{ \delta\mathbf{f}^T (\mathbf{f}^* - \mathbf{y}) \right\} > 0 \quad \forall \delta\mathbf{f} \neq \mathbf{0}$$

Dies wird erfüllt, wenn der zweite Term verschwindet, da der erste Term positiv definit ist.

Daraus resultiert eine notwendige Optimalitätsbedingung:

$$E \left\{ \delta\mathbf{f}^T (\mathbf{f}^* - \mathbf{y}) \right\} = 0$$

### Zwischenrechnung:

Mit

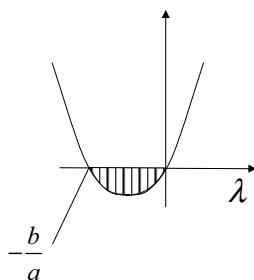
$$a := E \left\{ \|\delta\mathbf{f}\|^2 \right\} \quad \text{und} \quad b := 2E \left\{ \delta\mathbf{f}^T (\mathbf{f}^* - \mathbf{y}) \right\}$$

ergibt sich unter der Annahme einer Variation in eine beliebige Raumrichtung von  $\delta\mathbf{f}$ :

$$\delta\mathbf{f} \rightarrow \lambda\delta\mathbf{f}$$

$$a\lambda^2 + b\lambda > 0 \quad \forall \lambda$$

$$\lambda(a\lambda + b) > 0 \quad \forall \lambda$$



Ist  $b \neq 0$ , so gilt für alle  $\lambda \in (-\frac{b}{a}, 0)$ :

$$\lambda(a\lambda + b) < 0$$

$\Rightarrow b$  muss Null sein, dann lautet die Parabel:

$$a\lambda^2 > 0$$

Der Erwartungswert kann mit der Verbundverteilungsdichte ausformuliert werden zu:

$$E\{\delta \mathbf{f}^T(\mathbf{f}^* - \mathbf{y})\} = \sum_{\mathbf{x}} \sum_{\mathbf{y}} \delta \mathbf{f}^T(\mathbf{f}^* - \mathbf{y}) \cdot p(\mathbf{x}, \mathbf{y})$$

$$= \sum_{\mathbf{x}} \delta \mathbf{f}^T \left[ \sum_{\mathbf{y}} (\mathbf{f}^* - \mathbf{y}) \cdot p(\mathbf{y} | \mathbf{x}) \right] p(\mathbf{x}) = \mathbf{0}$$

Dieser Ausdruck verschwindet, falls der Term in der eckigen Klammer Null wird.

Daraus resultiert die folgende Optimalitätsbedingung:

$$\sum_{\mathbf{y}} (\mathbf{f}^* - \mathbf{y}) \cdot p(\mathbf{y} | \mathbf{x}) = \mathbf{f}^* \sum_{\mathbf{y}} p(\mathbf{y} | \mathbf{x}) - \sum_{\mathbf{y}} \mathbf{y} \cdot p(\mathbf{y} | \mathbf{x})$$

$$= \mathbf{f}^* - \sum_{\mathbf{y}} \mathbf{y} \cdot p(\mathbf{y} | \mathbf{x}) = \mathbf{0}$$

Die optimale Schätzfunktion ist die sogenannte Regressionsfunktion:

$$\mathbf{f}^*(\mathbf{x}) = \sum_{\mathbf{y}} \mathbf{y} \cdot p(\mathbf{y} | \mathbf{x}) = E\{\mathbf{y} | \mathbf{x}\}$$

## Polynomiale Regression

Wahl von Polynomen als Basisfunktionen für die Funktionsapproximation.

Zunächst für eine skalare Funktion mit einem vektoriellen Argument:

$$f(\mathbf{x}) = a_0 + a_1 x_1 + a_2 x_2 + \dots + a_N x_N$$

$$+ a_{N+1} x_1^2 + a_{N+2} x_1 x_2 + a_{N+3} x_1 x_3 + \dots$$

$$+ a_{\dots} x_1^3 + a_{\dots} x_1^2 x_2 + a_{\dots} x_1^2 x_3 + \dots$$

Dieses verallgemeinertes Polynom enthält eine Konstante  $a_0$ , gefolgt von  $N$  linearen Termen,  $N(N+1)/2$  quadratischen Termen usw.

Ein Polynom vom Grade  $G$  und der Dimension  $N$  des Argumentvektors  $\mathbf{x}$  hat

$$L = \binom{N+G}{G} = \frac{(N+G)!}{G!N!}$$

polynomiale Terme, welche als Basisfunktionen  $f_i(\mathbf{x})$ ,  $i=1,2,\dots,L$  einer Funktionsentwicklung betrachtet werden können.

Obiges Polynom kann kompakt durch Einführung einer vektoriellen Abbildung eines  $N$ -dimensionalen Vektors  $\mathbf{x}$  auf einen  $L$ -dimensionalen Vektor  $\mathbf{p}$  beschrieben werden:

$$\mathbf{p}(\mathbf{x}) = \left[ 1 \quad x_1 \quad x_2 \quad \dots \quad x_N \quad x_1^2 \quad x_1 x_2 \quad x_1 x_3 \quad \dots \quad x_1^3 \quad x_1^2 x_2 \quad x_1^2 x_3 \quad \dots \right]^T$$

Wir führen einen Koeffizientenvektor  $\mathbf{a}$  ein und erhalten für das skalare Polynom:

$$f(\mathbf{x}) = \mathbf{a}^T \mathbf{p}(\mathbf{x}) = \sum_{i=0}^L a_i p_i(\mathbf{x})$$

Für die Musterklassifikation benötigen wir für jede Klasse eine Polynomfunktion

$$f_k(\mathbf{x}) = \mathbf{a}_k^T \mathbf{p}(\mathbf{x}) \quad \text{verantwortlich für die Klassen} \\ k = 1, 2, \dots, K$$

Kombiniert man die klassenspezifischen Koeffizientenvektoren  $\mathbf{a}_k$  zu einer Matrix

$$\mathbf{A} = [\mathbf{a}_1 \quad \mathbf{a}_2 \quad \dots \quad \mathbf{a}_K]$$

so erhält man für die vektorielle Polynomfunktion:

$$\mathbf{f}(\mathbf{x}) = \mathbf{A}^T \mathbf{p}(\mathbf{x})$$

## Optimale Anpassung der Matrix $\mathbf{A}$

Es gilt nun, die Koeffizientenmatrix  $\mathbf{A}$  während der Optimierungsprozedur bestmöglich anzupassen.

Dabei kann das Approximationsproblem direkt auf die Messungen oder aber auf die Merkmale eines Unterraumes (z.B. nach einer KLT) angewendet werden.

Adaption der Koeffizientenmatrix  $\mathbf{A}$ :

$$J = \min_{\mathbf{f}(\mathbf{x})} E \left\{ \|\mathbf{f}(\mathbf{x}) - \mathbf{y}\|^2 \right\} \approx \min_{\mathbf{A}} E \left\{ \|\mathbf{A}^T \mathbf{p}(\mathbf{x}) - \mathbf{y}\|^2 \right\}$$

Das Variationsproblem mit der Suche nach einer unbekanntem Funktion wird durch die Wahl polynomialer Basisfunktionen auf ein Parameteroptimierungsproblem reduziert!

Das Gütekriterium  $J$  minimiert die Varianz des Residuums, nämlich den Fehler zwischen  $\mathbf{y}$  und der Approximation  $\mathbf{A}^T \mathbf{p}(\mathbf{x})$ .

Die Lösung wird bestimmt unter Ausnutzung des Variationsansatzes, was zu einer notwendigen Bedingung führt:

$$J(\mathbf{A}^* + \delta \mathbf{A}) \geq J(\mathbf{A}^*) \quad \text{für } \forall \delta \mathbf{A} \neq \mathbf{0}$$

Es gilt nun, die Koeffizientenmatrix  $\mathbf{A}$  während der Optimierungsprozedur bestmöglich anzupassen.

Adaption der Koeffizientenmatrix  $\mathbf{A}$ :

Wegen

$$J(\mathbf{A}) = E \left\{ \left\| \mathbf{A}^T \mathbf{p}(\mathbf{x}) - \mathbf{y} \right\|^2 \right\} = E \left\{ \left[ \mathbf{A}^T \mathbf{p} - \mathbf{y} \right]^T \left[ \mathbf{A}^T \mathbf{p} - \mathbf{y} \right] \right\}$$

$$\text{und} \quad \underbrace{\mathbf{a}^T \mathbf{b}}_{\text{Skalarprodukt}} = \text{Spur}(\underbrace{\mathbf{b} \mathbf{a}^T}_{\substack{\text{dyad.} \\ \text{Produkt}}})$$

$$\text{sowie} \quad \text{Spur}(\mathbf{PQ}) = \text{Spur}(\mathbf{QP})$$

erhält man:

$$\begin{aligned} J(\mathbf{A}) &= E \left\{ \text{Spur} \left[ \left[ \mathbf{A}^T \mathbf{p} - \mathbf{y} \right] \left[ \mathbf{A}^T \mathbf{p} - \mathbf{y} \right]^T \right] \right\} \\ &= \text{Spur} \left[ E \left\{ \mathbf{y} \mathbf{y}^T \right\} \right] - 2 \text{Spur} \left[ \mathbf{A}^T E \left\{ \mathbf{p} \mathbf{y}^T \right\} \right] + \text{Spur} \left[ \mathbf{A}^T E \left\{ \mathbf{p} \mathbf{p}^T \right\} \mathbf{A} \right] \\ &= E \left\{ \left\| \mathbf{y} \right\|^2 \right\} - 2 \text{Spur} \left[ \mathbf{A}^T E \left\{ \mathbf{p} \mathbf{y}^T \right\} \right] + \text{Spur} \left[ \mathbf{A}^T E \left\{ \mathbf{p} \mathbf{p}^T \right\} \mathbf{A} \right] \end{aligned}$$

Wie zuvor angenommen, kann der Zielvektor o.B.d.A. als Einheitsvektor formuliert sein und damit gilt:

$$J(\mathbf{A}) = 1 - 2 \text{Spur} \left[ \mathbf{A}^T E \left\{ \mathbf{p} \mathbf{y}^T \right\} \right] + \text{Spur} \left[ \mathbf{A}^T E \left\{ \mathbf{p} \mathbf{p}^T \right\} \mathbf{A} \right]$$

Mit dem Variationsansatz ergibt sich:

$$\begin{aligned} J(\mathbf{A}^* + \delta \mathbf{A}) &= 1 - 2 \text{Spur} \left[ (\mathbf{A}^* + \delta \mathbf{A})^T E \left\{ \mathbf{p} \mathbf{y}^T \right\} \right] \\ &\quad + \text{Spur} \left[ (\mathbf{A}^* + \delta \mathbf{A})^T E \left\{ \mathbf{p} \mathbf{p}^T \right\} (\mathbf{A}^* + \delta \mathbf{A}) \right] \\ &= 1 - 2 \text{Spur} \left[ \mathbf{A}^{*T} E \left\{ \mathbf{p} \mathbf{y}^T \right\} \right] - 2 \text{Spur} \left[ \delta \mathbf{A}^T E \left\{ \mathbf{p} \mathbf{y}^T \right\} \right] \\ &\quad + \text{Spur} \left[ \mathbf{A}^{*T} E \left\{ \mathbf{p} \mathbf{p}^T \right\} \mathbf{A}^* \right] + 2 \text{Spur} \left[ \delta \mathbf{A}^T E \left\{ \mathbf{p} \mathbf{p}^T \right\} \mathbf{A}^* \right] \\ &\quad + \text{Spur} \left[ \delta \mathbf{A}^T E \left\{ \mathbf{p} \mathbf{p}^T \right\} \delta \mathbf{A} \right] \end{aligned}$$

Eingesetzt in die Ungleichung der notwendigen Bedingung für ein Optimum ergibt:

$$\text{Spur} \left[ \delta \mathbf{A}^T E \left\{ \mathbf{p} \mathbf{p}^T \right\} \delta \mathbf{A} \right] - 2 \text{Spur} \left[ \delta \mathbf{A}^T (E \left\{ \mathbf{p} \mathbf{y}^T \right\} - E \left\{ \mathbf{p} \mathbf{p}^T \right\} \mathbf{A}^*) \right] \geq 0$$



Da der erste Term positiv definit ist und wegen der Positiv-Definitheit von der Momentenmatrix  $E\{\mathbf{p}\mathbf{p}^T\}$ , muss der zweite Term verschwinden:

$$\text{Spur}\left[\delta\mathbf{A}^T(E\{\mathbf{p}\mathbf{y}^T\}-E\{\mathbf{p}\mathbf{p}^T\}\mathbf{A}^*)\right]=0 \quad \forall \delta\mathbf{A}$$

und daraus schließlich die in  $\mathbf{A}$  lineare Gleichung:

$$E\{\mathbf{p}\mathbf{p}^T\}\mathbf{A}^* = E\{\mathbf{p}\mathbf{y}^T\}$$

diese Gleichung garantiert ein minimales Residuum des Fehlervektors:

$$\Delta\mathbf{f}(\mathbf{x}) = \mathbf{f}(\mathbf{x}) - \mathbf{y} = \mathbf{A}^T\mathbf{p}(\mathbf{x}) - \mathbf{y}$$

Die Lösung des Problems führt auf die Aufgabe, die beiden Momentenmatrizen  $E\{\mathbf{p}\mathbf{p}^T\}$  und  $E\{\mathbf{p}\mathbf{y}^T\}$  aus dem Trainingsdatensatz zu schätzen und eine lineare Matrixgleichung zu lösen. Der verbleibende Fehlervektor berechnet sich zu:

$$J(\mathbf{A}^*) = 1 - \text{Spur}\left[\mathbf{A}^{*T}E\{\mathbf{p}\mathbf{p}^T\}\mathbf{A}^*\right]$$

## Orthogonalität des Schätzfehlervektors

Die Bestapproximation fordert nach dem Projektionssatz, dass der Fehlervektor  $\Delta\mathbf{f}$  senkrecht steht auf den Unterraum, welcher durch die polynomialen Basisfunktionen in  $\mathbf{p}(\mathbf{x})$  aufgespannt wird.

In statistischer Notation bedeutet das, dass die Momentenmatrix, welche aus  $\mathbf{p}$  und  $\Delta\mathbf{f}$  aufgespannt wird, eine Nullmatrix sein muss:

$$E\{\mathbf{p}\Delta\mathbf{f}^T\} = E\left\{\mathbf{p}\left(\mathbf{A}^{*T}\mathbf{p} - \mathbf{y}\right)^T\right\} = E\{\mathbf{p}\mathbf{p}^T\}\mathbf{A}^* - E\{\mathbf{p}\mathbf{y}^T\} = 0$$

Diese Eigenschaft wird auch von  $\mathbf{p}$  auch auf  $\mathbf{f}$  vererbt, da  $\mathbf{f}$  aus Linearkombinationen der Basisfunktionen in  $\mathbf{p}(\mathbf{x})$  aufgebaut wird.

$$E\{\mathbf{f}^* \Delta\mathbf{f}^T\} = E\left\{\mathbf{A}^{*T}\mathbf{p}\Delta\mathbf{f}^T\right\} = \mathbf{A}^{*T}E\{\mathbf{p}\Delta\mathbf{f}^T\} = 0$$

## Erwartungstreue der Approximationsfunktion

Die Schätzfunktion ist erwartungstreu (unbiased). Der Schätzfehler (Residuum) hat den Erwartungswert  $\mathbf{0}$ .

$$E \{ \Delta \mathbf{f} \} = E \{ \mathbf{f}(\mathbf{x}) - \mathbf{y} \} = \mathbf{0}$$

Daraus folgt, dass  $\mathbf{f}(\mathbf{x})$  eine erwartungstreue Schätzung von  $\mathbf{y}$  ist:

$$E \{ \mathbf{f}(\mathbf{x}) \} = E \{ \mathbf{y} \}$$

## Rekursive Lernregel für den Polynomklassifikator

Für den auf das minimale Fehlerquadrat aufbauende Klassifikator ergibt sich die folgende rekursive Lernstrategie für die Koeffizientenmatrix  $\mathbf{A}$ :

$$\mathbf{A}_n = \mathbf{A}_{n-1} - \alpha \left[ \sum_{i=1}^n \mathbf{p}_i \mathbf{p}_i^T \right]^{-1} \mathbf{p}_n \left[ \mathbf{A}_{n-1}^T \mathbf{p}_n - \mathbf{y}_n \right]^T$$

Diese Lernregel beinhaltet folgende Schritte:

- Wahl eines Startwerts für  $\mathbf{A}$
- Die momentane Stichprobe  $[\mathbf{x}, \mathbf{y}]$  wird abgebildet mit Hilfe der Polynombasis  $\mathbf{p}(\mathbf{x})$  auf das Paar  $[\mathbf{p}, \mathbf{y}]$
- Berechne die Schätzung  $\mathbf{f}$  auf der Basis der gegebenen Beobachtung  $\mathbf{p}$  und der momentanen Koeffizientenmatrix  $\mathbf{A}$ :  $\mathbf{f} = \mathbf{A}^T \mathbf{p}$
- Berechne das Residuum:  $\Delta \mathbf{f} = \mathbf{f} - \mathbf{y}$
- Korrigiere  $\mathbf{A}_n$  auf der Grundlage obiger Rekursion

$$\mathbf{A}_n = \mathbf{A}_{n-1} - \alpha \mathbf{G}^{-1} \mathbf{p}_n \Delta \mathbf{f}^T$$

Die Gewichtsmatrix  $\mathbf{G}$  ist dabei die Momentenmatrix  $E\{\mathbf{p}\mathbf{p}^T\}$  auf der Basis der  $n$  zur Verfügung stehenden Stichproben.

Eine strikte Beachtung obiger Rekursion, erfordert eine weitere Rekursion für die Inverse der Gewichtsmatrix gemäß:

$$\mathbf{G}_n^{-1} = \frac{1}{1-\alpha} \left[ \mathbf{G}_{n-1}^{-1} - \alpha \frac{\mathbf{G}_{n-1}^{-1} \mathbf{p}_n \mathbf{p}_n^T \mathbf{G}_{n-1}^{-1}}{1 + \alpha (\mathbf{p}_n^T \mathbf{G}_{n-1}^{-1} \mathbf{p}_n - 1)} \right]$$

$\mathbf{G}$  oder die echte zusätzliche Iteration ändern fast nichts am Ergebnis.  $\mathbf{G}$  kann sogar noch vereinfacht werden zu:

$$\mathbf{G} = E\{\mathbf{p}\mathbf{p}^T\} = \sum_{i=1}^n \mathbf{p}_i \mathbf{p}_i^T$$

## Herleitung der rekursiven Lernregel

$$E\{\mathbf{p}\mathbf{p}^T\} \mathbf{A} = E\{\mathbf{p}\mathbf{y}^T\}$$

Einführung von Schätzwerten:

$$\left( \underbrace{\frac{1}{n} \sum_{i=1}^n \mathbf{p}_i \mathbf{p}_i^T}_{\mathbf{G}_n} \right) \mathbf{A}_n = \underbrace{\frac{1}{n} \sum_{i=1}^n \mathbf{p}_i \mathbf{y}_i^T}_{\mathbf{H}_n}$$

$$\mathbf{G}_n \mathbf{A}_n = \mathbf{H}_n \quad \text{und} \quad \begin{aligned} \mathbf{G}_n &= (1-\alpha) \mathbf{G}_{n-1} + \alpha \mathbf{p}_n \mathbf{p}_n^T \\ \mathbf{H}_n &= (1-\alpha) \mathbf{H}_{n-1} + \alpha \mathbf{p}_n \mathbf{y}_n^T \end{aligned}$$

$$\Rightarrow \mathbf{G}_n \mathbf{A}_n = (1-\alpha) \mathbf{H}_{n-1} + \alpha \mathbf{p}_n \mathbf{y}_n^T$$

$$= (1-\alpha) \mathbf{G}_{n-1} \mathbf{A}_{n-1} + \alpha \mathbf{p}_n \mathbf{y}_n^T$$

$$= (\mathbf{G}_n - \alpha \mathbf{p}_n \mathbf{p}_n^T) \mathbf{A}_{n-1} + \alpha \mathbf{p}_n \mathbf{y}_n^T$$

$$= \mathbf{G}_n \mathbf{A}_{n-1} - \alpha \mathbf{p}_n (\mathbf{p}_n^T \mathbf{A}_{n-1} - \mathbf{y}_n^T)$$

$$\Rightarrow \boxed{\mathbf{A}_n = \mathbf{A}_{n-1} - \alpha \mathbf{G}_n^{-1} \mathbf{p}_n \left[ \mathbf{A}_{n-1}^T \mathbf{p}_n - \mathbf{y}_n^T \right]^T}$$

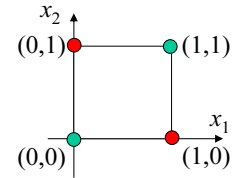
# Polynomklassifikator für das XOR-Problem

Wir wählen als Regressionsfunktion ein Polynom mit Termen zweiten Grades in der Hoffnung, dass damit eine Lösung erreicht werden kann:

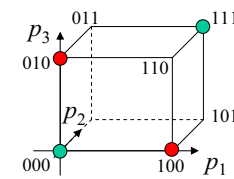
$$f(\mathbf{x}) = a_1x_1 + a_2x_2 + a_3x_1x_2 = \mathbf{a}^T \mathbf{p}(\mathbf{x})$$

mit:

$$\mathbf{p} = [x_1, x_2, x_1x_2]^T \quad \text{und} \quad \mathbf{a} = [a_1, a_2, a_3]^T$$



$\mathcal{X}$ -Raum



$\mathcal{P}$ -Raum

Die Eckpunkte des Quadrats im zwei-dimensionalen Ursprungsraum  $\mathcal{X}$  werden auf die Ecken eines (Hyper-)Würfels einen dreidimensionalen Merkmalsraum  $\mathcal{P}$  abgebildet und dieser schließlich auf den eindimensionalen Bedeutungsraum  $\mathcal{Y}$ :

$$\mathcal{X} \rightarrow \mathcal{P} \rightarrow \mathcal{Y}$$

$\mathcal{X}$	$x_1$	0	1	0	1
	$x_2$	0	0	1	1
$\mathcal{P}$	$p_1$	0	1	0	1
	$p_2$	0	0	1	1
	$p_3$	0	0	0	1
$\mathcal{Y}$	$y$	0	1	1	0

Die optimalen Parameter der Regressionsfunktion  $\mathbf{a}^*$  ergeben sich aus:

$$\underbrace{E\{\mathbf{p}\mathbf{p}^T\}}_{\mathbf{R}_{pp}} \mathbf{a}^* = \underbrace{E\{\mathbf{p}\mathbf{y}\}}_{\mathbf{R}_{py}} \quad \text{bzw.:} \quad \mathbf{a}^* = \mathbf{R}_{pp}^{-1} \cdot \mathbf{R}_{py}$$

$$\text{mit: } \mathbf{R}_{pp} = E\{\mathbf{p}\mathbf{p}^T\} \approx \frac{1}{4} \sum_{i=1}^4 \mathbf{p}_i \mathbf{p}_i^T$$

$$= \frac{1}{4} \left\{ \mathbf{0} + \begin{bmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} + \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 0 \\ 1 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix} + \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} \right\} = \frac{1}{4} \begin{bmatrix} 2 & 1 & 1 \\ 1 & 2 & 1 \\ 1 & 1 & 1 \end{bmatrix}$$

$$\text{und: } \mathbf{R}_{py} = E\{\mathbf{p}\mathbf{y}\} \approx \frac{1}{4} \sum_{i=1}^4 \mathbf{p}_i y_i = \frac{1}{4} \sum_{i=1}^4 y_i \mathbf{p}_i$$

$$= \frac{1}{4} \left\{ 0 \cdot \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} + 1 \cdot \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} + 1 \cdot \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} + 0 \cdot \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \right\} = \frac{1}{4} \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}$$

Daraus folgt:

$$\mathbf{a}^* = \mathbf{R}_{pp}^{-1} \cdot \mathbf{R}_{py} = \begin{bmatrix} 1 & 0 & -1 \\ 0 & 1 & -1 \\ -1 & -1 & 3 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ -2 \end{bmatrix}$$

und damit eine Regressionsfunktion

$$f(\mathbf{x}) = a_1 x_1 + a_2 x_2 + a_3 x_1 x_2 = x_1 + x_2 - 2x_1 x_2$$

welche an den 4 Funktionswerten genau die Werte der Zielfunktion annimmt und ansonsten interpoliert!