

```
1 Entscheidungstheorie
2 =====
3
4
5 1. Wkeit
6 -----
7
8 * bedingte
9
10 * totale
11
12 * Satz von Bayes
13
14 * a priori
15
16 * a posteriori
17
18
19 2. Bayesche Entscheidungsregel
20 -----
21
22 * stetige Featurewerte
23   - bedingte Wkeitdichte
24   - Dichte
25
26
27 * Featurevektor
28
29
30 3. Fehlerrate
31 -----
32
33 * Formel
34
35 * Bayessche
36
37
38 4. Risiko-Minimierung
39 -----
40
41 * Wertung von Features
42
43 * Entsch.Regel zur RM
44
45
46 5. Diskriminantenfunktion
47 -----
48
49 * Vorteile
50   - Vereinfachung Entscheidungsregel
51   - Konzept für disk. Ansatz
52
53
54 * Template Matching / Rocchio Klassifiz.
55
56 * lineare
57
58
59 6. mehrdim/multivariate NormVert.
60 -----
61
62 * Kovarianzmatrix
63   - Formel
64   - Aus Training: empirische Kov.
65
66 * Wdichte
67
68 * Varianz aller Features gleich: lin. DF ist Template-Matching.
69 Idee: ln() nutzen und Konst. bzw von i unabh. Summanden verwerfen
70
71
72 Wkeit-Dichte
73 =====
74
75 1. nicht parametrische Dichteschätzung
76 -----
```

```
77
78
79 * Suche Schätzer für p(x) ohne Annahme zur Verteilung.
80 Also: kein Schätzer für Parameter von p(x)
81
82 1. Idee
83   - Histogramm
84   - Anzahl Messungen in Intervall
85   - Bei Binominalverteilung: erw.treuer Schätzer für p in Intervall
86   - Probleme
87   - d Dimensionen: N^d Intervalle
88   - viele Zellen mögl. leer
89   - Lösungen
90   - unstetig? Kernel-Methoden
91     - Zu viele oder leere Zellen? variable Zellengröße
92
93
94 * Kernel-Methoden
95   - keine willkürlichen Intervallgrenzen
96   - Wkeit eines Messwerts wird aus Distanz innerhalb einer Fensterbreite
97   - rund um alle Messwerte ermittelt
98   - gleitendes Histogramm möglich:
99     p an jedem Punkt x statt nur in einem Intervall
100   - mit n Messw. und K - positiv und auf 1 norm. - ist das auch
101     für Dichte pDach(x)
102   - pDach(x) = 1/(n*h) SumAlle.n( K((x-xi)/h) )
103   - d-Dim: kugelsym. Kern mit (2pi)^(-d/2) * e^(-<x,x>/2)
104
105
106 * Fensterbreite h
107   - Kernel: globales h
108     - MeanSquareError = Bias^2(pDach(x)) + Var(pDach(x))
109     - Güte Schätzung über alle durch Integration: MISE(pDach)
110   - TyLorentw.: von MISE mit pDach(x) in h:
111     - AMISE(h) Ableitung = 0 für Minimum in h
112     - Bei NV: h = 0.79 * IQR * n^(-1/5)
113     - Bei NV, weil h sonst zu sehr gedrückt: Silvermans Rule of Thumb
114     - h_srot = 0.9 * min{sigmaDach, IQR/1.35} * n^(-1/5)
115   - Plugin Methode: numerische Lösung von Minimum h aus
116     - AMISE'(h) = 0 mit h_srot als Startpunkt
117
118
119 * KNN-Methoden
120   - statt variabler Anz. Messwerte in festem h:
121     - feste Anzahl der Messwerte in var. h
122   - lokales h: Vermeidung, dass h bei hoher Wkeit zu klein/niedrig
123   - statt h jetzt: V(x) = 2* (k-kleinsten Abstand von x)
124
125   - Fehlerrate
126   - Erwartungswert der Fehlerrate asymptotisch Abschätzbar
127     - e B = Bayessche Fehlerrate
128     - Cover & Hart: blödsster Fall k=1 aber n -> inf liefert:
129       e B < E(e) < 2e B
130
131   - pDach(x) = k/(n*sum x - xj l) leider keine Wkeit-Dichte!
132   - Integral +- inf ist nicht 1
133   - mit k-1 statt k lässt sich aber mit steigendem n zeigen, dass
134     fDach n(x) sich der WDichte in x nähert
135   - k(n) ist Nachbarzahl
136   - fDach n(x) = (k(n)-1)/(n*V(x))
137   - Nicht so tragisch! Wir wollen x ja einer Klasse aus
138     Trainingsdaten zuordnen
139   - für a posteriori Wkeit P(wi | x) reicht's
140   - KNN-Entsch.Regel: x wird der Class zugeteilt, die unter
141     den k-nächsten Nachb am häufigsten ist!
142   - gleich häufig: 1) zufällig, 2) am nächsten an x, 3) im Mittel
143     am Nächsten
144
145   - Parameterwahl
146     - k nie größer als niedrigste "Klassenpopulation"
147     - Enas & Choi: k ca n^(2/8)
148     - für gutes k: Gütekriterien nutzen
149     - Wahl der MEIRIK: Was heißt "Nächster"? Ggf. Bewertung von
150       Features untersch.
151
152   - Edit Trainingsdaten
153   - Ziele
```

```
153 - Verringern Fehlerrate
154 - Beschl. Klassifizierung
155
156 - Methoden
157
158 - einzelne Ausreißer entfernen
159 - Ausreißer werden als fehlklassifiziert angesehen, sobald
160 sie nicht mit dem Rest via KNN-Regel klassifizierbar sind
161
162 - Condensing
163 - entfernen von überflüssigen Samples
164 - Entfernung darf bei restlichen Samples die "Decision Boundaries"
165 nicht ändern
166 - Algorithm oft von Reihenfolge der Samples abh.
167 - Ab 3 Class: Minimierung in konsistente Teilmenge NP-vollst.
168
169 - Prototype-Generator
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228

- Maximum Likelihood Schätzer
- wie ist Parameter Theta zu wählen, so dass aus beob. Werte
  max. Wkeit haben
- Da jede Messung unabh.: Maximiere Produkt aus Wkeitverteilung
  je beobachteten Wert -> theta Dach
- Gradient(Nabla_v) von Likelihood-Fkn = 0
- Vorsicht: Wendepkt? Rand von Theta?
- mit log-Likelihood-Fkn ergibt 01.2 bei NV: Mittelwert und emp. Var.
- Mittelwert-Vektor und (nicht erw.treue) Kovar.Matr

* erw. Treue Schätzer
- Erw.Wert des Schätzers sollte Parameter liefern
- BIAs: Abweichung Erwartungswert des empirischen Schätzers von
  theor. Schätzers
- SigmaDach = (1/(n-1)) * SummeAlle.k von
  (vek.x k - vek.muDach)*(vek.x k - vek.muDach)^t

* Bayes Classifier NV
- Nutzung emp. Kovarianz und Mittelwert
- naivesBayes(): Wkeit x wenn in wi ist Produkt der Wkeit je Feature,
  wenn Feature in wi (=Feature unabh.)
- Listing gaussian.bayes():
  Nutzung der Trainingsdaten für Wkeiten von x k
- Aber: jede Klasse ist hier gleich-Wscheinlich
- a priori: aus Erfahrung bzw. Anfangs-Wkeiten

- Alternativ: Kerneldichte-Schätzer für Wkeiten nutzen:
  NICHT PARAMETRISCH!

3. Quantil
-----
* IQR: Inter Quartile Range x_0.75 - x_0.25
* x_0.25 ist der x-Wert, bei dem 25% der Messwerte xi kleiner sind
* Quartil: 0.75 oberes, 0.5 mittleres (Median), 0.25 unteres

Classifier
=====
1. diskriminatorischer Ansatz
-----
* class ist Teilregion des Featurespace
```

```
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304

2. probabilistischer Ansatz
-----
* class und feature ist P zugeordnet
* allg. als disk. Ansatz

3. Features
-----
* Testdaten
* Trainingsdaten
* Extraktion

4. Modell auf Tr.Data: Verallgemeinerung
-----
5. Modell schwer bei Overfitting
-----
6. Vergleich von 2 (A und B)
-----
* Naiv
- Fehlerrate bei GLEICHEN Testdaten messen
- Unterschied der Fehlerate könnte auch an zufälliger Verteilung
  der Testdaten liegen
- je Testdaten könnte Unterschied auch bei 2 identischen
  Classifizieren auftreten

* Nutzung Signifikanzniveau: ähnlich Konfidenzniveau

* Testdaten unab. zuf. gewählt
- A(xi) und B(xi) unabhängig
- A(xi) und B(xi) NICHT unabhängig, weil beide mit selben
  Testsample gemessen werden müssen
- McNemar Test
- Fehleratenvergleich
- Betrachtet werden nur die Anzahl der Fehlklassifizierungen,
  wo A richtig und B falsch (n01) liegt bzw. umgekehrt (n10)
- Wenn beide gleich gut, sollte in etwa gelten:
  n10 = n01 = (n10 + n01) / 2
- Shi2 = Maß des Unterschieds der Fehler
  = ( | n10 - n01 | -1 )^2 / (n10+n01)
- Shi(m)^2 (1-alpha) = 3.8415 mit Signifikanzniveau alpha=0.05,
  dem (1 - 0.05)-Quantil und m=1 als Anz. der Freiheitsgrade
- Beispiel: ist Shi2 > 3.8415, dann ist Wkeit für rein statische
  Abweichung der Fehlerraten < 5%

- Paired t-Test
- Vergleich stetiger Performancemaße (nicht eindeutig Richtig/Falsch)
- Differenz D der Performace je Sample i wird betrachtet:
  Di = A(xi_i) - B(xi_i)
- Betrachte Teststatistik zu Zuf.Variable D:
  t=(Dquer * wurz(n) )/s(D)
  - Dquer = Mittelwert aller i=1..n Di
  - s^2(D) = (n-1)^-1 * Sum (Di -Dquer) ^2
  - Gleiche Performance wäre: E(D)=0

- signifikant untersch. Performance, wenn Betrag von
  T > t(n-1) (1-alpha/2)
- t(n-1) (1-alpha/2) ist das (1-alpha/2) Quantil der t-Verteilung
  mit n-1 Freiheitsgraden
- Mit steigenden Freiheitsgrad (mehr Messwerte) nähert sich
  Quantil dem der Std.NV
- Quantil der Std. NV ist immer kleiner oder gleich

7. Receiver Operation Characteristic (ROC)
-----
- 4 -
```

```
-----
* Vergleich von Classifier, die Samples nur 2 Klassen zuordnen
* Classifier unterscheiden sich von nur 1 Parameter (Z.B. k in KNN
  oder ein Schwellwert für Entscheidung)
* ROC-Kurve
  - jeder Punkt stammt von einem anderen Parameter
  - AUC: Area Under Curve:
    je größer die Fläche, desto besser der Classifier
  - x-Achse: epsilon2 = Rate der Falsch als Klasse w2 Zugeordneten
  - y-Achse: 1-epsilon1 ist Rate, der korrr. w1 Zugeordneten
* bester Parameter?
  - Nutzung von Kosten bei Fehlklassifizierung
  - geringste Kosten (Minimum) sind da, wo Steigung P(w2)/P(w1) ist
  - Kostengrade am weitesten von "teurer" Diagonale entfernt, aber
    tangiert trotzdem noch ROC-Kurve
Metrische Räume
=====
1. Metrik d auf M bildet M x M -> R ab
-----
* d(x,y) = 0 <=> x = y
* Symetrie: Abstand von x nach y ist gleich wie von y nach x
* Dreiecksungleichung: Ein Umweg könnte länger sein!
2. falls Menge Vektorraum
-----
* V x V -> R ist Skalarprodukt, wenn
  - bilinear: <x+a,y> = <a,y> + <a,y> und <c x,y> = c<x,y>
  - symmetrisch: <x,y> = <y,x>
  - positiv: <x,x> > 0 (wenn x!=0)
* Skalarprodukt kann Norm definieren: ||x|| = wurzel(<x,x>)
* eine Norm kann eine Metrik definieren: d(x,y) = || x-y ||
* allg. R^n Skalarprodukt: <x,y> = vek.x.Transponiert * A * vek.y
3. Metrik Beispiele
-----
* triviale Metrik: d(x,y) ist 0 wenn x=y sonst 1
* Euklidische Metrik: d(x,y) = sqrt(sum(xi - yi)^2)
* Levensthein Distance: d(x,y) = minimale Anz. Edit-OP um aus
  x das Wort y zu machen
4. Norm
-----
* Abbildung eines Elements aus Vektorraum V über R
* N1: ||x|| = lcl * ||x||
* N2: ||x+y|| <= ||x|| + ||y||
* N3: ||x|| = 0 <=> x ist Nullvektor
* Lp-Norm oder Minkowski-Norm
  - euklid p=2
  - Manhattan p=1
  - Maximumsnorm: p -> inf
-----
```

```
-----
- p-te Wurzel aus Summe der Beträge aller (xi)^p
5. Parallelogramm-Gleichung
-----
* || x+y ||^2 + || x-y ||^2 = 2 ||x||^2 + 2 ||y||^2
* dann gilt: Norm stammt von Skalarprodukt
* <x,y> = 1/4 ( || x+y ||^2 - || x-y ||^2 )
6. Nutzung Freiheit des SkProd. (A) für Anpassung Metrik
-----
* Diagonalmatrix mit >=0 Werten -> Feature Weighting
* cm oder Meter: Feature Normalization
* feste Varianz: je Feature auf 1 normiert wi = 1/(si)^p
  - A=Inv(sigma) Inverse der Kovarianzmatrix bei unab. Features
  - Metrix: MAHALANOBIS Distanz
* Normierung auf Mittelwerte = 0 ist nur Verschiebung:
  beeinflusst nicht Abstand
HMM
===
1. allgemein
-----
* durchläuft Folge von Class
* Sprache: Folge von Phonemen
* Schrift: t = Schreibrichtung
* Gebärde: Folge von Handpos.
* Wkeit für Folgen berechnen: Markov-Ketten
2. Begriffe HMM
-----
* wi = Hidden States
* vi (i=1..S) Visible States
* aij = Wkeiten der Hidden States
* bjk = Wkeit, dass im State wj der Wert vk emittiert wird
3. HMM Evaluation
-----
* Anwendung
  - aij,bjk, Startzustand w(0) gegeben
  - Wkeit für generierung der Folge: v(1)...v(T)
* Naiv: leider 0(T * C^T)
* FORWARD-Algorithmus
  - Mit Startwerten wird Wkeit berechnet für v(1) für C viele
    Hidden States
  - iterativ: für t+1 wird wieder je C aus (Summe der bisherigen
    Wkeiten der Hidden States)* bjk die Wkeit für v(t+1) berechnet.
  - Wegen Nutzung vorheriger Berechnungen: O(T*C^2)
4. aij und bjk ? (j,k Laufindex für C,S)
-----
```

```
-----
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532

5. HMM Dekodierung
-----
* Anwendung
  - aij, bjk gegeben sowie eine Folge v(1)..v(T)
  - Was wahrscheinlichste Folge w(1)..w(T) die zu Beobachtung v führt?

* Naiv: Brute-Force alle w(1)..w(T) durchprobieren O(C^T)

* VITERBI-Algorithmus
  - Prinzip ähnlich wie Forward-Algorithmus
  - statt Aufsummierung der Wkeiten wird nun bisher Max Übergangswkeiten
    abgespeichert (Produkt max(letzte Wkeit)*Übergang jetzt)
  - am Ende via Backtracking die Indizes der Classes des zuletzt Max.
    Pfades auflisten
  - Wegen Nutzung vorheriger Berechnungen: O(T*C^2)

6. HMM Training/Learning Problem
-----
* Anwendung
  - Anzahl(=C) der w und Anzahl(=S) der v
  - Auch gegeben: Trainingsdaten mit mehreren beobachteten Folgen v(1)..v(T)
  - Wahrsch. Parameter aij, bjk ?
  - P(w(1) = wi) also Wkeit der HiddenStates zu begin ??

* Naiv: Maximum-Likelihood und Parameterschätzung
  - liefert keine lösbar analytische Formel
  - liefert keinen bekannten algorithmischen Lösungsweg

* Algorithmus
  - Numerisch/Iterativ
  - Irgendwelche Anfangswerte nehmen
  - Dekodierung der Trainingsdaten (Ermittlung
    Wahrsch. Hidden-state-Folge)
  - Schätze Parameter neu:
    aij=(Anz. Übergänge wi -> wj)/(Anzahl Vorkommen wi)
  - Wiederhole Dek. bis irgendein Konvergenzkrit. erfüllt ist

7. in der PRAXIS
-----
* left-right Modell
  - wi -> wj nur, wenn j>= i
  - Wenn Buchstabe "e" mehrfach im Wort, dann bekommt er eigenen
    Hidden State

* Isolated Word Recognition
  - Jedes Wort hat ein HMM(aij,bjk,pi)
  - Jedes Wort wird trainiert (aij,bjk)
  - pi = Wkeit des Hidden State wi
  - unbekanntes Wort: Evaluation, welches HMM am Wahrscheinlichsten
    unbekanntes v emittiert. Also: max( P(v, wenn aij,bjk gilt) )
  - -> Hidden Markov Toolkit (HTK)

* ein
  - 4 Hidden States
  - aij: Matrix nicht ergodisch, wenn nur obere Dreiecksmatrix
    (für jedes A^n bleibt untere Hälfte =0)
  - Beispiel Zeile 2 der 4x4 Matrix: 0 a22 a23 0
  - Erste 0: es passiert bei dem Wort nicht, dass nach "i" wieder
    zurück auf "e" gewechselt wird
  - a22: Wkeit, dass "i" tangezogen wird: eiiin_
  - a23: Wkeit, dass nach "i" ein "n" kommt:
  - letztes 0: Wkeit, dass man das "n" bei "ein" vergisst
```

```
-----
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608

Fehlerabschätzung / Classifier Performance Estimation
=====
1. Allgemein
-----
* Classifier: grossPhi(x) ordnet Feature X einer Class zu
* Baysche/Optimale Fehlerrate: e B = e ergibt sich, sollte Classifier
  immer die Wahrscheinlichste Klasse wählen
* Wahrsch. Klasse unbekannt, weil Wkeit-Verteilung in Praxis unbekannt
* Fehlerrate e = E( F_grossPhi(X,w)
* F_grossPhi(X,w) = 0 falls grossPhi(x) = w , sonst 1 (falsch klassif.)
* tatsächliche/true Errorrate e T
  - will man abschätzen
  - ist Maß, wie gut Trainingsdaten in der Lage sind, Features
    zu klassifizieren
  - expected Error Rate e_E: Mitteln über alle Trainingssets der Größe
    n ergibt einen Erwartungswert für e_T

2. Holdout Estimate
-----
* Herausheben von Trainingsamples: ES GIBT TESTDATEN
* bilden mehrerer Trainingsamples der Größe m
* eDach T = 1/m * SummeAlleSamples( F_grossPhi(xi,wi) )

3. Konfidenzniveau / Irrtumswkeit
-----
* Abweichung Schätzer und "true" soll unwarscheinlicher sein als Alpha
* P( l p - pDach l >epsilon) <= alpha
* "true" p liegt mit Wkeit (1-alpha) in Intervall pDach +/- epsilon
  - asym: Agresti-Coull
  - asym: HPD

* KLASSISCH: via Normalverteilung
  - pDach = k/m
  - epsilon = z (1 - alpha/2) * sigma
  - sigma^2 = 1/m * p(1-p)
  - Trick: nutzte für p in sigma das pDach !
  - Wegen Trick: sigma als Breite ist falsch!

* AGRESTI-COULL KonfidenzIntervall
  - bei alpha = 5%
  - m und pDach werden durch "quantil" modifiziert
  - m -> m+4
  - p ist nicht mehr k/m sondern: (k+2)/(m+4)

* Highest Posterior Density HPD
  - ebenfalls Abschätzung von e_T
  - Wkeit für error p hängt nicht nur von k sondern selbst von p ab
  - a = k+1
  - b = m-k +1
  - Nutzung Beta-Verteilung Be(p,a,b)
  - um mit alpha an p1,p2 zu kommen:
  - Bestimmungsgleichungen lösen
    - Be(p1,a,b) = Be(p2,a,b)
    - p1.Integral.p2 Be(p,a,b) dp = 1-alpha

4. ohne Testdaten
-----
```

\* zu wenig manuell klassifizierte Daten?!

\* Nutzen aller Daten zum Training und/oder Testen

\* Resubstitution / Apperent Error Rate

- Trainingsdaten = Testdaten
- zu optimistisch
- Entscheidungsgrenzen an Trainingsdaten angepasst
- GANZ ÜBEL bei KNN mit  $k=1$ :  $e = A=0$

\* Cross-Validation (Leave-One-Out)

- allgemein ist k-fold: aus Trainingsdaten werden k Teilmengen gebildet
- k-fold:  $n=k \rightarrow$  In jeder Teilmenge wurde nur ein Wert weggelassen
- e CV = Mittelwert aus allen Fehlklassifizierungen

\* Bootstrap Estimate

- Schätzung der Varianz des Schätzers
- Schätzung des Bias von e\_A
- Algorithmus
- K-oft aus Trainingsdaten T n Samples ziehen mit zurücklegen  $\rightarrow$  Sk
- Je Sk: trainiere Classifier
- Je Sk: e Ak bilden mit Sk+gebildeten Classifiere
- Je Sk: e Tk bilden mit Sk+Classifiere, der aus T gewonnen wurde
- Je Sk: Bk = e Tk - e Ak
- eBach boot =  $e + 1/K * \text{SummeAlle Bk}$

-  $B = E(e_T - e_A) \leq e_T = e_A + B$

\* Andere Bootstrap-Methoden

- unterschied, wie Bias "gemittelt" wird
- Erfahrung: Linearkombination untersch. e\_A
- unterschied, wie Trainingsdaten "gezogen" werden

Markov-Ketten  
=====

1. Stochastische Matrix  
-----

\* Einträge im Intervall  $[0,1]$

\* Summe der Zeileneinträge = 1

\* die i-te Zeile: Spalte j gibt Wkeit an, von i nach j zu wechseln

2. Stochastischer Vek.  
-----

\* Summe der Einträge ist =1

3. Wkeit  $\pi_i(t) = P(w_i(t)) = \text{SummeAlle}(\pi_j(t-1) * a_{ji})$   
-----

\* Also: wegen Unabhängigkeit ist es Summe Wkeit vorheriger Zustand \* "Übergangs-Wkeit"

\* Iterativ:  $\pi_i(t)$  lässt sich aus AnfangsWkeit  $\text{vek}.p(0) * A^t$  (a-Matrix t-mal angewand) berechnen

\* Statt  $\pi_i(t) \rightarrow \text{vek}.p(t)$  STOCHASTISCHER VEKTOR

\* Statt  $a_{ji} \rightarrow$  A STOCH. MATRIX

\*  $\text{vek}.p(t) \rightarrow$  Zeilenvektor

4. Übergangsmatrix  
-----

\* allgemein könnte  $a_{ij}(t)$  gelten.

\* homogene/stationäre MK:  $a_{ij}$  nicht von t abhängig

5. markovsche Eigenschaft  
-----

\* Wkeit eines Zustands zum Zeitpunkt t hängt nur von Wkeit des Zustands  $w(t-1)$  ab

\* w älter als t-1 ist nicht relevant (Gedächtnislosigkeit)

\* EINE Markov-Kette ist daher mit einer Übergangsmatrix komplett beschreibbar. Jeder Eintrag ist Wkeit für Zustandswechsel

6. stochastischer Prozess  
-----

\* Abfolge von Zuständen: ZUFÄLLIG, nicht von Zustandsfolge abhängig

\* zur Zeit t ist System in Zustand  $w_i(t)$

\* System muss zur Zeit t in einem Zustand sein:  
Summe über alle  $w_i(t) = 1$

7. Allgemein  
-----

\* System durchläuft Class

\* Class = Zustand des Systems

\* zu diskre. Zeiten ändert sich (zufällig) Zustand (oder bleibt gleich)

8. Langzeitverhalten  
-----

\* ist A strikt positiv, gilt:  
- mit jeder  $\text{vek}.p(0)$  und  $t \rightarrow \infty$ :  $p(0) * A^t \rightarrow p$   
- p hat dann Eigenschaft:  $p * A = p$

\* nur ein  $A^k$  muss für ein k strikt Positiv sein: A ist ERGODISCH

Rejection & Konfidenzmaße  
=====

1. Zurückweisen  
-----

2. Ambiguity Rejection  
-----

\* selbst die Class, wo Sample am ehesten zu passt, ist nicht sehr Wahrscheinlich

\* Sample liegt auch nahe einer anderen Class

\* oft, wenn mit einzelnen Feature nicht gut die Klassen unterscheidbar sind

\* Finden eines Konfidenzwerts für Feature  
-  $\text{fi}(x) < f_0$  dann wird zurückgewiesen, wenn Schwellwert  $f_0$  nicht erfüllt

- Finden eines Schwellwerts: z.B. mit Leave-One-Out, bis Errorrate mit Trainingsdaten minimal ist

\* Chow's Reject Option  
- geht davon aus, das a posteriori Wkeit  $P(w_i|x)$  bekannt ist  
- Fehlerrate immer  $e(t) \leq t$   
- Reject nur, wenn  $t < 1 - 1/C$

```
761 - Satz von Chow
762 - Zusammenhang Threshold, Fehlerrate, Akzetanzrate
763 - x wird zurückgewiesen, wenn  $\text{MAX.i}(P(w_i, x)) < 1 - t$ 
764 -  $r(t) = \text{Reject-Rate}$ 
765 -  $e(t) = 0. \text{Integral.t } r(s) \text{ ds} - t * r(t)$ 
766
767 - NAIV: optimales t wenn Accuracy A(t) maximal
768 - max bei  $t=0$ 
769 - als Optim.Kriterium dumm
770 -  $t=0 \rightarrow$  alles wird zurückgewiesen
771
772 - OPTIMAL Reject Threshold: QUOTIENT aus Kosten für ein Reject und
773 Kosten für ein Fehler
774
775
776 3. Distance Rejection
777 -----
778
779 * Sample so weit von restlichen Trainingsdaten entfernt, dass vermutlich
780 zu keiner Class passt
781
782 * Novelty Detection: Ist Sample eine neue, eigene Class?
783
784 * Nutzung einer Distanz ist wie ein "Konfidenzwert", aber hier sind
785 es fiktive Kanten/Regionen statt "echte" Nachbarn
786 -  $g(x) = 1/k * \text{SummeAlleNachbarnY } d(x,y)$ 
787
788
789 * Bestimmung der empirischen Verteilung, wie Wahrscheinlich es ist,
790 dass Sample den Abstand d zum Nachbar hat
791 -  $F(d) = P(g(x) \leq d)$ 
792 - z.B. mit Trainingsdaten via LeaveOneOut
793 - Es lässt sich durchschnittliche max. Distanz d max bestimmen
794 -  $1 - F(d \text{ max}) \leq \epsilon$ 
795 - epsilon = Anteil der Daten, die zurückgewiesen werden
796 - epsilon weist aber auch korrekte Daten zurück und bestimmt so
797 False-Positive Rate
798
799
800 4. A(t) = P(correct | accept)
801 -----
802
803 * = P(correct und accept) / P(accept)
804
805
806 * = (  $1 - e(t) - r(t)$  ) / (  $1 - r(t)$  )
807
808 Feature Selection & Extraction
809 =====
810
811 1. Ziele
812 -----
813
814 * Verringerung Error Rate
815
816 * Dimensionsreduzierung
817
818 * Selection: Teilmenge nehmen
819
820 * Extraction: Neukombination
821
822
823 2. 1) Wrapper-Methoden
824 -----
825
826 * Classifier-Performance-Schätzer
827
828 * z.B. Leave-One-Out
829
830 * zu Aufwändig wegen Training (außer bei kNN)
831
832 3. 2) Filter-Methoden: Maß unabh. von Classifier
833 -----
834
835 * gute Samples
836 - grosser Abstand between
```

```
837 - kleiner Abstand within
838
839
840 * Ges. Maß für Güte der Class
841 - Mittelwert 2er Class
842 - Min 2er Class
843
844
845 * Class-Separation via Abstandsmaß
846 - 1) Theoretisch
847 - Chernoff Schranke beschränkt e B nach oben
848 - Spezialfall: NV, und  $s = 0.5$ 
849 - Bhattacharyya Distanz B als Abstandsmaß von w1 und w2
850 -  $e_{CB} = B * \text{wurzel}( P(w1) * P(w2) )$ 
851 - für B muss man Mittelwerte und Kovarianzmatrizen aus
852 Trainingsdaten schätzen
853
854 - 2) Training: Streumatrix
855 -  $S = n * K_{\text{dach}} = (n-1) * \text{SigmaDach}$ 
856 - S zerlegen in Streuung between Class und innerhalb (within)
857 -  $SW = \text{SumAlleClass}[ ni * K_{\text{dach.i}} ]$ 
858 -  $SB = \text{SumAlleClass}[ n.i * (\mu_i - \mu)(\mu_i - \mu)^t ]$ 
859 - leider nur gut bei unimodal Verteilungen - stark um Mittelwert
860 - J1 = Spur(SB)/Spur(SW)
861 - J2 = Spur(  $SW^{-1} * SB$  )
862
863
864 4. Algorithmen zur Selection
865 -----
866
867 * Ziel: Welche Teilmenge von f Features optimiert J?
868
869 * Brute Force:  $2^f$ 
870
871 * max d viele Features:  $O(f^d)$ 
872
873 * Skalar
874 - Idee: Performance nur in einem Feature Messen
875 und die Besten Features nehmen
876 - AUC der ROC-Kurve
877 - Fishers Discriminant Ratio
878 - gut, wenn man 2 Dim auf eine Abbilden will
879 -  $FDR = (\mu_1 - \mu_2)^2 / (\sigma_1^2 + \sigma_2^2)$ 
880
881 - obides: nur für 2 Klassen. Trick: Vergleich via Mittelwerte oder
882 Min der ermittelten J
883
884
885 * Vektor
886 - Greedy
887 - Gefahr, im ersten Extremum hängen zu bleiben
888 - Seq Backward Sel.: Lasse immer mehr Features weg, jenachdem
889 welches Weglassen die Güte am meisten erhöht
890 - Seq Forw. Selection: nehme immer mehr Features dazu
891
892 - Monte-Carlo Methoden
893 - springt gelegentlich auf Extrema heraus
894 - Features werden wie Chromosomen gemischt
895
896
897 5. Extraction mit Principal Component Analysis (PCA)
898 -----
899
900 * Hauptachsentransformation
901
902 * Idee: Suche im Mehrdim. Featurespace Mittelpunkt und Profeziere
903 Punkte auf Grade, die durch Mittelpkt geht.
904
905 * Abstand der Punkte zu der Grade sollte um ganzen Minimal sein
906
907 * Die Hauptachse einer Kovarianzmatrix ist in Richtung des Eigenvektors,
908 der zum größten Eigenwert gehört. Die Daten Streuen entlang dieser
909 Richtung am Stärksten!
910
911 * Independent Component Analysis (ICA) nutzt diese Richtungen, um ohne
912 Klasseninfo Richtung der Stärken Abweichung von NormalVerteilung
```

```
913 zu finden. -> Ist es dann eine neue Klasse?
914
915
916 6. Lineare Diskriminanten Analyse (LDA)
917 -----
918
919 * im Gegensatz zu PCA: Nutzung von Klasseninfo! 2 Klassen sind Thema
920
921 *  $J(\text{vek.w}) = \text{FDR} = \frac{\text{abs}(m1 - m2)^2}{(s1^2 + s2^2)}$ 
922
923 *  $m_i = \langle \text{vek.w}, \text{vek.mi} \rangle$ 
924
925 *  $m_i$  = Schwerpkt Koordinaten, die auf vek.w projiziert sind
926
927 *  $\text{abs}(m1 - m2)$  = Abstand der Klassmittltpkt proj. auf vek.w
928
929 *  $s_i^2 = \text{emp. Var. der Class } i \text{ längst vek.w}$ 
930
931 * Gesucht w1 und w2 (vek.w)
932
933 * vek.w ist in Richtung  $\text{SW}^{-1} * (\text{vek.m1} - \text{vek.m2})$ 
934
935
936 Clustering
937 =====
938
939 1. Unsupervised Learning
940 -----
941
942 2. Visual. hochdim. Daten
943 -----
944
945 * PCA (Hauptachsentransf.) mit Euklidischer Metrik
946
947 * Multidim. Scaling (MDS)
948   - Lösungsweg via Stress-Fktn
949   - Sammons Nonlinear Mapping
950   - start mit PCA
951   - gehe zu nächsten lok. Minimum entlang des stärksten Gradienten
952
953   - isoMDS von Kruskal
954
955   - MDS ist ein  $f(\text{vek.x})$ , so dass:
956      $d(\text{vek.xi}, \text{vek.xj}) \approx d(f(\text{vek.xi}), f(\text{vek.xj}))$ 
957
958
959 3. Finden von Klassen
960 -----
961
962 * flach
963   - kmeans
964     - Fail: durchprobieren alle Möglichen K Teilmengen
965     -  $J(K) = \text{Spur}(\text{SW})$ 
966     - Minimum von  $J(K)$ 
967     - leider NP-hart
968
969   - Isodata Algorithmus
970     - findet nur lokale Minima, könnte aber mit Monte Carlo
971       Verfahren verbessert werden
972     - es müssen Start-Mittelpkte gewählt werden
973     - Mittelpkt-Vektor konvergiert
974     - Es werden immer wieder Teilmengen gebildet, wo enthaltene
975       x näher an einem Mittelpkt sind, als an Mittelpunkten
976       anderer Teilmengen
977     - Mittelpkte der Teilmengen werden nach jeder Zusammenstellung
978       von Teilmengen neu berechnet
979
980   - Wahl der Clusterzahl K bei d Dimensionen und n Samples
981     - a) Ellebogen: Stärkster Knick in  $J(K)$ -Kurve
982     - b) info.theor.Methode
983       - minimum von:  $-2 \cdot l(K) + a(n) \cdot d * K$ 
984       -  $l(K)$  = Maximum der log-Likelihood Funktion
985       -  $d * K$  = Anzahl Parameter des Modells
986     - a(n) = Strmaß für Komplexität des theor. Models
987       - Bayesian Information Criterion ( $\text{BIC}$ ) =  $\ln(n)$ 
988       - Akaike Information Criterion ( $\text{AIC}$ ) = 2
```

```
989
990
991
992   - c) Calinski-Harabasz Index Maximieren:
993      $\text{CH}(K) = \{ \text{Spur}(\text{SB}) / (K-1) \} / \{ \text{Spur}(\text{SW}) / (n-K) \}$ 
994
995   - graph-basiert: MinSpannTree
996     - Menge zerfällt in Teilmengen durch entfernen Inkonsistenter Kanten
997     - Kante e inkonsistent:  $\text{kanntengewicht} > m_{ü.e} + q * \text{sigma.e}$ 
998     -  $q = 1.96 \dots$  ist Analogie zu 95% Konfidenzintervall
999     -  $m_{ü.e}$  und  $\text{sigma.e}$ : Wird aus Kanten berechnet, die max k Schritte
1000       von e Entfernt sind
1001     - K muss hier nicht vorher festgelegt werden!
1002
1003 * hierarchisch
1004   - agglomerieren mit Dendrogram
1005   - Zusammenlegen von x, die sich am nächsten sind
1006   -  $\text{cdist} = \text{Cluster-Distanz}$ 
1007   - gewünschte Clusterzahl K: wenn n-K oft Daten/Cluster
1008     zusammengefasst wurden
1009
1010   - single/complete Link Clustering
1011     - kürzester Abstand der Cluster wird aus kürzestem Abstand
1012       zwischen den Clustern genommen
1013     - kürzester Abstand der Cluster wird aus WEITESTEM Abstand
1014       zwischen den Clustern genommen
1015
```