

Skript zur Vorlesung  
**Knowledge Discovery in Databases II**  
im Wintersemester 2012/2013

**Lecture 2: Feature Selection und Feature Reduktion**

Lectures: Dr Matthias Schubert, Dr. Eirini Ntoutsis  
Tutorials: Erich Schubert

Notes © 2012 Matthias Schubert,

[http://www.dbs.ifi.lmu.de/cms/Knowledge\\_Discovery\\_in\\_Databases\\_II\\_\(KDD\\_II\)](http://www.dbs.ifi.lmu.de/cms/Knowledge_Discovery_in_Databases_II_(KDD_II))

**2. Featurereduktion und Featureselektion**

Inhalt dieses Kapitels

2.1 Einführung

Motivation, „*Curse of Dimensionality*“

2.2 Feature-Selektion

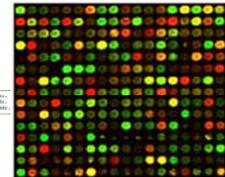
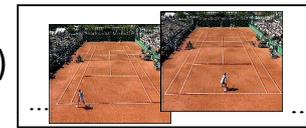
Methoden zur Auswahl geeigneter Unterräumen

2.3 Feature-Reduktion

Transformation der Feature-Räume

### Beispiele für hochdimensionale Objektdarstellungen

- Bilddaten
  - Bilddeskriptoren (Farbverteilungen, Pyramid of Gradients, Texturen,..)
  - Aufteilung der Bilder in Bereiche
  - Je nach Typ ca. 16 - 1000 Features
- Metabolomdaten
  - Feature entspricht Konzentration von Stoffwechselprodukt im Blut
  - Je nach Aufbau 50-2000 Features
- Mikro-Arrays
  - Features entsprechen z.B. Genen
  - Je nach Versuchsaufbau bis zu 20000 Features
- Texte
  - Features entsprechen Worten oder Termen
  - Je nach Anwendung 10000-20000 Features

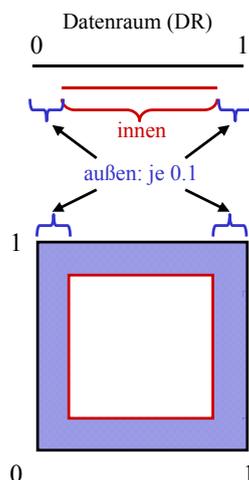


```

1 Kennzeichnung: Name:         Anzahl:
2 Teamraum:                   Menge:  1
3                                     Platz:  1
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
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
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
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
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
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
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
    
```

### „Fluch der Dimensionalität“ (Curse of Dimensionality)

- Distanz zum nächsten Nachbarn unterscheidet sich im Hochdimensionalen kaum von der Distanz zu einem beliebigen Nachbarn
- Die Wahrscheinlichkeit, dass Datenobjekte am Rand des Datenraumes liegen, steigt mit der Anzahl der Dimensionen exponentiell



1D:  $P_{DR} = 1^1 = 1$   
 $P_{innen} = 0.8^1 = 0.8$   
 $P_{außen} = 1 - 0.8 = 0.2$

2D:  $P_{DR} = 1^2 = 1$   
 $P_{innen} = 0.8^2 = 0.64$   
 $P_{außen} = 1 - 0.64 = 0.36$

3D:  $P_{DR} = 1^3 = 1$   
 $P_{innen} = 0.8^3 = 0.512$   
 $P_{außen} = 1 - 0.512 = 0.488$

10D:  $P_{außen} = 0.893$

andere Effekte des *Curse of Dimensionality*:

- Die Werte in jeder Dimension sind verrauscht.  
D.h. die Werte schwanken durch Störeinflüsse, die mit dem Objekt an sich nichts zu tun haben.
- Bei zunehmender Dimensionalität wird die Summe der Störeinflüsse so groß, dass die beobachteten Unterschiede zwischen 2 Objekten von den Störeinflüssen dominiert werden.  
=> Summe der Unterschiede hängt von der Ausprägung der Störeinflüsse ab und nicht von den Objekteigenschaften.  
=> Abstand zwischen Objekt gleicht sich immer weiter an, da Störeinflüsse gleichverteilt über alle Objekte.

- Patterns und Modelle auf hochdimensionalen Daten sind oft schwer interpretierbar.  
⇒ Lange Entscheidungsregeln
- Effizienz bei hochdimensionalen Daten häufig problematisch.  
⇒ Indexstrukturen degenerieren auf hochdimensionalen Daten  
⇒ Distanzberechnungen werden i.d.R. teurer
- Muster treten nur in Teilräumen auf, aber nicht im gesamten Featureerraum
- Cliques von korrelierten Features dominieren das Objektverhalten.

**Idee:** Bei sehr vielen Features sind nicht alle notwendig, um die Daten zu beschreiben:

- Features sind nicht aussagekräftig für ein Problem
- Informationsgehalt stark korrelierter Features ist fast identisch

Die Einschränkung auf einen Teilraum des gesamten Datenraums kann Verfahren effizienter und effektiver machen.

**Lösung:**

Streiche alle „überflüssigen“ Dimensionen aus dem Feature-Raum.

**Gegeben:** Vektorraum  $F = D_1 \times \dots \times D_n$  mit  $D = \{D_1, \dots, D_n\}$ .

**Gesucht:** Minimaler Unterraum  $M$  über  $D \subseteq D$ , der für ein gegebenes Data Mining Problem eine optimale Lösung erlaubt.

=> Minimalität erhöht Effizienz und verringert den Curse-of-Dimensionality.

Probleme:

- Optimalität hängt unmittelbar mit dem anschließend verwendeten Data Mining Algorithmus zusammen.
  - Problem ist sehr komplex, da es  $2^d$  mögliche Unterräume gibt.
  - Die meisten Qualitätsmaße haben keine Monotonie-Eigenschaften.  
(Features können erst in Kombination mit anderen Features nützlich werden)
- => Suche läuft meist über Heuristiken, die für ein gegebenes Qualitätsmaß nicht den besten sondern nur einen guten Raum findet

1. Forward Selektion und Feature Ranking  
Information Gain ,  $\chi^2$ -Statistik, Mutual Information
2. Backward Elimination und Zufällige Unterräume  
Nächste Nachbarn Kriterium, Modellbasierte Suche
3. Branch and Bound Suche mit Inkonsistenzmaß
4. Genetische Algorithmen zur Unterraum-Suche
5. Feature Clusterings für Unsupervised Probleme

**Gegeben:** ein Supervised Learning Task

- eine Zielvariable  $C$  soll vorhergesagt werden
- es gibt Trainingsobjekte für die Relation zwischen Feature-Raum und Klassenvariable

**Vorgehen**

- Berechne Güte  $G(D_i, C)$  für jede Dimension  $D_i \in \{D_1, \dots, D_n\}$  und die Zielvariable  $C$
- Sortiere Dimensionen  $\{D_1, \dots, D_n\}$  nach Güte  $G(D_i, C)$
- Wähle die  $k$  besten Dimensionen

**Annahme:**

Feature sind außer über die Zielvariabel nicht direkt korreliert  
=> es reicht aus die Korrelation zwischen Feature und Zielvariable zu untersuchen

Wie gut kann die Zielvariable bei gegebenen Feature-Wert vorhergesagt werden?

Bewertung der Features über statistische Maße:

- Basieren auf Verteilungen bzgl. der Klassen und der Feature-Werte  
Achtung bei reellen Features ist Schätzung einer Zufallsvariable nicht direkt durchführbar. Daher:
  - Split zur Umwandlung von reellen Features in diskrete Features
  - Annahme über Verteilungsfunktion (z.B. Normalverteilung,..)
- Wie stark korrelieren die Verteilungen von Klassen und Features?
- Wie stark unterscheidet sich Aufteilung bzgl. dieses Features von einer zufälligen Aufteilung ?

**Idee:** Bewerte wie gut jede Dimension die Klassen „unterscheidet“.  
(vgl. Entscheidungsbäume)

Zerlege Trainingsmenge anhand Feature/Unterraum in Teilmengen  
(Unterteilung: nach Werten oder Splitkriterien).

Die *Entropie* für eine Menge  $T$  von Trainingsobjekten ist definiert als

$$entropie(T) = -\sum_{i=1}^k p_i \cdot \log p_i \quad (p_i \text{ steht für Häufigkeit der Klasse } i \text{ in } T)$$

$$entropie(T) = 0, \text{ falls } p_i = 1 \text{ für ein } i$$

$$entropie(T) = 1 \text{ für } k = 2 \text{ Klassen mit } p_i = 1/2$$

$$\text{Informationsgewinn}(T, a) = entropie(T) - \sum_{j=1}^m \frac{|T(a)_j|}{|T|} \cdot entropie(T(a)_j)$$

mit  $T(a)_j$  Teilmenge von  $T$ , für die Attribut  $a$  den Wert  $j$ -ten Wert annimmt.

Bei reell wertigen Attributen muss ein Split gefunden werden.

**Idee:** Bewertet Unabhängigkeit einer Dimension von einer Klasse.

Aufteilung der Daten anhand der Splitwerte  $s$  oder anhand diskreter Attributwerte

$$A = \left| \{o \mid x \leq s \wedge \text{Class}(o) = C_j\} \right| \quad \text{Objekte in } C_j \text{ mit Wert } x \leq \text{Splitwert}$$

$$B = \left| \bigcup_{l \neq j} \{o \mid x \leq s \wedge \text{Class}(o) = C_l\} \right| \quad \text{Objekte anderer Klassen mit } x \leq \text{Splitwert.}$$

$$C = \left| \{o \mid x > s \wedge \text{Class}(o) = C_j\} \right| \quad \text{Objekte in } C_j \text{ mit } x > \text{Splitwert.}$$

$$D = \left| \bigcup_{l \neq j} \{o \mid x > s \wedge \text{Class}(o) = C_l\} \right| \quad \text{Objekte anderer Klassen, mit } x > \text{Splitwert}$$

χ<sup>2</sup>-Statistik ist definiert durch: 
$$\chi^2(t, C_j) = \frac{|DB| (AD - CB)^2}{(A + C)(B + D)(A + B)(C + D)}$$

Je höher Maximum oder Durchschnitt über alle Klassen desto besser das Feature  $a$ :

$$\chi^2_{\max}(a) = \max_{i=1}^m \{ \chi^2(a, C_i) \} \quad \text{oder} \quad \chi^2_{\text{avg}}(a) = \sum_{i=1}^m \Pr(C_i) \chi^2(a, C_i)$$

Maß für gegenseitige Abhängigkeit zweier Zufallsvariablen.

**Hier:** Vergleich der Abhängigkeit von der allgemeinen Klassenverteilung mit Verteilung der Feature Vektoren.

1. diskreter Fall.

$$I(X, Y) = \sum_{y \in Y} \sum_{x \in X} p(x, y) \log \frac{p(x, y)}{p(x)p(y)}$$

2. Kontinuierlicher Fall

$$I(X, Y) = \int_Y \int_X p(x, y) \log \frac{p(x, y)}{p(x)p(y)} dx dy$$

Häufigkeit der Wertekombination  $x, y$

Term der die Abhängigkeit darstellt:  
bei Unabhängigkeit  $p(x, y) = p(x)p(y)$   
 $\Rightarrow \log(1) = 0$

### Vorteile:

- Effizient da nur einzelne  $d$  Features statt  $\binom{d}{k}$  Unterräumen untersucht werden
- Abschätzung der Verteilungen benötigt moderate Anzahl an Beispielen

### Nachteil:

- Dimensionen werden einzeln betrachtet: Klasse und Dimensionswert müssen direkt korreliert sein. Erkennt keine nicht achsenparallelen Muster.
- Korrelierte Dimensionen: Auswahl von bedeutungsgleichen Dimensionen falls diese am stärksten mit der Klasse korreliert sind

**Idee:** Starte mit dem gesamten Feature-Raum und lasse unnütze Feature weg.

**Vorgehen:** Greedy Backward Elimination

1. Generiere alle Unterräume  $R$  des Feature-Raums  $F$
2. Bewerte alle Unterräume  $R$  mit Gütekriterium für Räume  $G(R)$
3. Wähle den besten Unterraum  $R^*$
4. Fall  $R^*$  die Zieldimensionalität hat, wird Suche abgebrochen.  
ansonsten wird Algorithmus mit  $R^*$  aufgerufen.

**Anwendung:**

- kann für Supervized oder unsupervised Verfahren verwendet werden ( $G(R)$  kann sich auch an strukturellen Eigenschaften orientieren und braucht keine Zielvariable)
- Heuristische Suche, die nur Sinn ergibt wenn  $G(R)$  nicht monoton ist  
=> bei Monotonie kann exakt mit Branch and Bound gesucht werden.

**Idee:** Unterraum ist gut, wenn Abstand von Objekten innerhalb einer Klasse durchschnittlich kleiner ist, als zwischen Objekten unterschiedlicher Klassen.

**Qualitätsmaß:**

Für alle Objekte  $o \in DB$  berechne den nächsten Nachbarn in Klasse  $C=Class(o)$   $NN_C(o)$  und den kleinsten nächsten Nachbarn  $NN_{K \neq C}(o)$  in einer der anderen Klassen .

$$\text{Güte des Unterraums } U: Q(U) = \frac{1}{|DB|} \cdot \sum_{o \in DB} \frac{NN_{K \neq C}^U(o)}{NN_C^U(o)}$$

Bemerkung: Kriterium ist nicht monoton.

=> Durch Streichen einer Dimension kann Qualität besser oder auch schlechter werden.

**Idee:** Teste den Unterraum direkt damit, wie gut er für das gegebene Data Mining Problem funktioniert.

**Beispiel:** Bewerte jeden Teilraum mit einem Naive Bayes Klassifikator und 10-facher Überkreuzvalidierung auf einer Stichprobe von 200 Instanzen.

**Anwendung:**

- Erfolg des Data Mining Algorithmus muss messbar sein.  
(z.B. Klassifikationsgenauigkeit)
- Laufzeit sollte nicht zu hoch sein
- Parameter Auswahl sollte keinen zu hohen Einfluss haben
- Testmenge sollte überschaubar sein

### Vorteile:

- Betrachtet ganze Unterräume anstatt einzelner unabhängiger Dimensionen (mehrfache Abhängigkeiten der Zielvariable bleiben erhalten)
- Redundante Features können also solche erkannt und entfernt werden.

### Nachteile:

- Tests bzgl. der Qualität eines Unterraums sind meist wesentlich aufwendiger.
- Ansatz ist ebenfalls nur eine Heuristik da es keine Monotonie gibt.

**Gegeben:** Klassifikationsproblem über  $F$ .

**Ziel:** Selektiere  $k$  Dimensionen

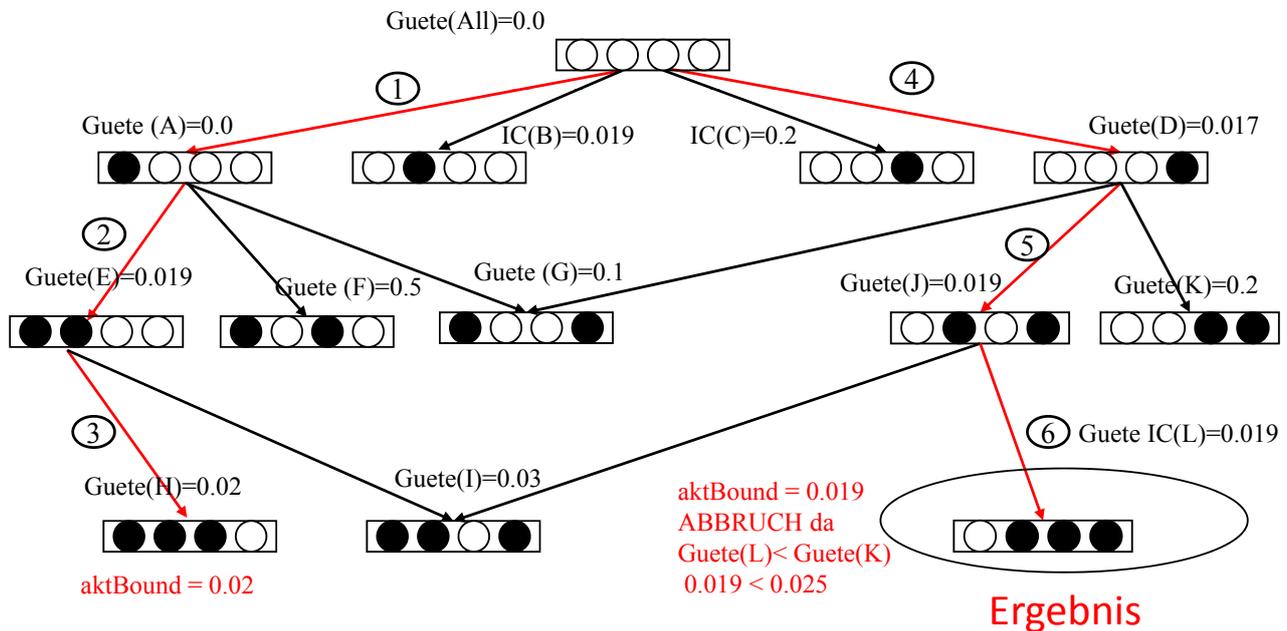
Backward-Elimination mit Branch and Bound:

```

FUNCTION BranchAndBound(Featurespace F, int k)
    queue.init(ASCENDING);
    queue.add(F, Guete(F))
    aktBound:= INFINITY;
    WHILE queue.NotEmpty() and aktBound > queue.top() DO
        aktURaum := queue.top();
        FOR ALL Unterräume U von aktURaum DO
            IF U.dimensionality() = k THEN
                IF Guete(U) < aktBound THEN
                    aktBound := Guete(U);
                    BestURaum := U;
            ELSE
                queue.add( U, Guete(U));
    RETURN BestURaum
    
```

Beispiel:  $k = 1$ .

○ Feature selektiert ● Feature ausgeschlossen



**Idee:** Existieren im Unterraum  $U$  identische Vektoren  $u, v$  mit  $v_i = u_i, 1 \leq i \leq d$  aber unterschiedlichen Klassen labels  $C(u) \neq C(v)$ .

=> Unterraum ist inkonsistent

Maß für die Konsistenz des Unterraums  $U$ :

- $X_U(A)$ : Anzahl aller zu  $A$  identischen Vektoren in  $U$
- $X_U^c(A)$ : Anzahl aller zu  $A$  identischen Vektoren in  $U$  die zur Klasse  $c$  gehören
- $IC_U(A)$ : Inkonsistenz bzgl.  $A$  in  $U$

$$IC_U(A) = X_U(A) - \max_{c \in C} X_U^c(A)$$

$$\text{Für ganz } U: IC(U) = \frac{\sum_{A \in DB} IC_U(A)}{|DB|}$$

Monotonie:  $U_1 \subset U_2 \Rightarrow IC(U_1) \geq IC(U_2)$

### Vorteile:

- Monotonie ermöglicht die effiziente Suche nach optimalen Unterräumen mit Branch and Bound
- Gut geeignet für Binäre Attribute, da hier die Wahrscheinlichkeit von identischen Vektoren hoch ist.

### Nachteil:

- Schlecht geeignet für numerische Attribute (ohne identische Vektoren nutzlos)
- Worst-Case Laufzeit immer noch exponentiell in  $d$

**Idee:** Wähle  $n$  zufällige Teilräume der Zieldimensionalität  $k$  aus  $\binom{d}{k}$  mögliche Teilräumen, teste diese und nimm den besten.

### Anwendung:

- Benötigt Qualitätsmaß für ganze Teilräume.
- Aufwand und Qualität hängen stark von der Anzahl der Sample-Teilräume  $n$  ab.

### Nachteil:

- Sucht nicht gezielt nach guter Kombination aussagekräftiger und unabhängiger Attribute
- Aufwand stark von der Wahl des Parameter  $n$  und angewendeten dem Qualitätsmaß abhängig.

**Idee:** Versuche gezielt hochqualitative Unterräume zu generieren.

**Ansatz:**

- **Population von Lösungen** := Menge  $k$ -dimensionaler Unterräume
- **Fitnesskriterium:** Qualitätsmaß für Unterräume
- **Mutationregel und Mutationwahrscheinlichkeit:**  
mit Wahrscheinlichkeit  $x\%$  wird Dimension  $D_i$  in  $U$  durch  $D_j$  ersetzt
- **Fortpflanzung:** Kombinationsregel für 2 Unterräume  $U1$  und  $U2$ :  
Wähle  $50\%$  der Dimensionen aus  $U1$  und  $50\%$  aus  $U2$
- **Selektionsregel:** Alle Kandidationräume die  $z\%$  schlechtere Fitness haben als der beste der bisherigen Generation sind nicht lebensfähig.
- **Freilos:** Zusätzlich zur Selektion kann jeder Unterraum mit Wahrscheinlichkeit  $u\%$  in die nächste Generation übernommen.

**Ablauf:**

Initialisiere Population

WHILE Max\_Fitness > Old\_Fitness DO

    Mutiere Population gemäß Mutationsrate

    WHILE nextGeneration < PopulationSize DO

        Generiere neuen Kandidaten  $K$  durch Fortpflanzung

        IF  $K$  hat Freilos oder  $K$  ist fit enough THEN

$K$  darf in die nächste Generation

    RETURN fittester Unterraum

### Bemerkung:

- hier nur Skizze des Grundalgorithmus (viele Erweiterungen)
- Konvergenz meist nur unter „Simulated Annealing“  
(Freiloswahrscheinlichkeit sinkt mit Anzahl der Generationen)

### Vorteil:

- Vermeidung von lokalen Maxima
- häufig gute Approximation des optimalen Unterraums

### Nachteile:

- kann lange Laufzeiten aufweisen
- viele Parameter müssen richtig gewählt werden, um einen guten Trade-Off zwischen Qualität und Laufzeit zu erzielen

**Gegeben:** Clusteringproblem über  $F$ .

**Ziel:** Reduziere Featurespace auf  $k$  Dimensionen.

**Vorgehen:** Da man irrelevante Attribute nicht erkennen kann, beschränkt man sich auf die Elimination von redundanter Information.

**Idee:** Clustere die Feature im Raum der Objekte und selektiere 1 Repräsentanten pro Cluster.

Maß für die Abhängigkeit Ähnlichkeit der Features:

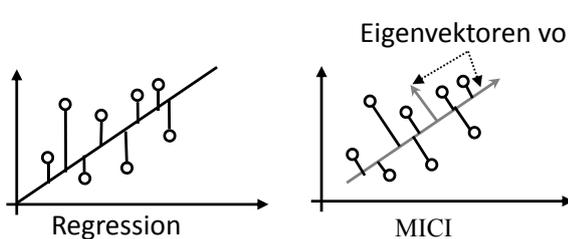
Korrelation zwischen 2 Features: 
$$COR(X, Y) = \frac{COV(X, Y)}{\sqrt{VAR(X)} \sqrt{VAR(Y)}}$$

Maximaler Information Compression Index (MICI):

**Idee:** Messe den kleinsten Eigenwert der Kovarianzmatrix  $\Sigma$  zwischen beiden verglichenen Unterräumen.

$$\begin{aligned} \det(\Sigma - \lambda E) &= \det \begin{pmatrix} \text{VAR}(X) - \lambda & \text{COV}(X, Y) \\ \text{COV}(X, Y) & \text{VAR}(Y) - \lambda \end{pmatrix} \\ &= (\text{VAR}(X) - \lambda) \cdot (\text{VAR}(Y) - \lambda) - \text{COV}(X, Y)^2 \\ \Rightarrow 0 &= \lambda^2 - \lambda \cdot \text{VAR}(Y) - \lambda \cdot \text{VAR}(X) + \text{VAR}(X) \cdot \text{VAR}(Y) - \text{COV}(X, Y)^2 \\ \Rightarrow \lambda &= \frac{(\text{VAR}(Y) + \text{VAR}(X)) \pm \sqrt{(\text{VAR}(Y) + \text{VAR}(X))^2 + 4 \cdot 1 \cdot \text{VAR}(X) \cdot \text{VAR}(Y) - \text{COV}(X, Y)^2}}{2 \cdot 1} \end{aligned}$$

$$\text{MICI}(X, Y) = \text{VAR}(Y) + \text{VAR}(X) - \sqrt{(\text{VAR}(Y) + \text{VAR}(X))^2 + 4 \cdot \text{VAR}(X) \cdot \text{VAR}(Y) - \text{COV}(X, Y)^2}$$



MICI:

- Symmetrisch
- Kann aus beiden Dimensionen 1 gebildet werden, die die gesamte Information beider widerspiegelt.

**Ablauf:**

- Cluster Feature mit  $k$ -medoid Clustering für  $k$  Cluster und Abstandsmaß MICI.
- Selektiere die Clusterrepräsentanten als Feature-Dimensionen

**Bemerkung:**

- Versucht für jede Gruppe abhängiger Dimensionen eine representative Dimension
- Anwendung anderer Clustering Algorithmen denkbar.  
(K-Means: Wähle Dimension die am nächsten am Cluster-Centroid liegt)
- Häufig werden Cluster-Algorithmen für Streams verwendet, wegen Ihrer Laufzeit  $O(d)$

### Vorteile:

- Verhältnismäßig schnelle Selektionsmethode
- Kommt ohne Klasseneinteilung aus

### Nachteile:

- Meist kein eindeutiges Ergebnis, da Clustering von Parametern und Reihenfolge abhängen kann.
- Repräsentative Dimensionen wechseln bei unterschiedlichen Cluster Algorithmen
- basiert auf paarweiser Korrelation  
=>höherwertige Dimensionen werden nicht untersucht.

- *Forward-Selection*: Untersuche einzelne Dimensionen  $D' \in \{D_1, \dots, D_d\}$ . und bilde neue Unterräume durch Kombination.
  - Greedy Selection mit Information Gain,  $\chi^2$  Statistik oder Mutual Information
- *Backward-Elimination*: Beginne mit dem gesamten Featurespace  $F$  und entferne überflüssige Dimensionen.
  - Greedy Elimination mit dem Modelbasierten oder Nächsten Nachbar Ansätzen
  - Branch and Bound Suche auf Basis des Inkonsistenzkriteriums
- *k-dimensionale Projektionen*: Suche direkt in der Menge der k-dimensionalen Unterräume nach einem geeigneten
  - Genetische Algorithmen zu Suche geeigneter Teilräume (Qualität wie beim Backward Elimination)
  - Feature Clustering mit dem MICI

### Diskussion:

- Viele Algorithmen basierend auf unterschiedlichen Heuristiken
- Feature können aus 2 Gründen eliminiert werden:
  - es existieren andere bedeutungsgleiche Feature (Redundanz)
  - Features sind nicht mit der Aufgabe korreliert
- häufig können auch schon nicht optimale Ergebnisse sowohl Effizienz als auch Effektivität verbessern
- **Vorsicht:** Selektierte Features müssen keinen direkten Einfluß auf Zielvariable haben, sondern können auch nur von den gleichen versteckten Einflüssen abhängen.

- Gründe für Feature Selektion
- Curse-of-Dimensionality
- Forward Selection und Feature Ranking
  - Information Gain
  - $\chi^2$  Statistik
  - Mutual Information
- Backward Elimination
  - Modelbasierte Güte
  - Nächste Nachbar Methode
  - Inkonsistenz und Branch & Bound
- k-dimensionale Projektionen
  - Genetische Algorithmen
  - Feature Clustering

- A. Blum and P. Langley: *Selection of Relevant Features and Examples in Machine Learning*, Artificial Intelligence (97), 1997
- H. Liu and L. Yu: *Feature Selection for Data Mining* (WWW), 2002
- L.C. Molina, L. Belanche, Â. Nebot: *Feature Selection Algorithms: A Survey and Experimental Evaluations*, ICDM 2002, Maebashi City, Japan
- P. Mitra, C.A. Murthy and S.K. Pal: *Unsupervised Feature Selection using Feature Similarity*, IEEE Transactions on pattern analysis and Machine intelligence, Vol. 24. No. 3, 2004
- J. Dy, C. Brodley: *Feature Selection for Unsupervised Learning*, Journal of Machine Learning Research 5, 2004
- I. Guyon, A. Elisseeff: *An Introduction to Variable and Feature Selection*, Journal of Machine Learning Research 3, 2003
- M. Dash, H. Liu, H. Motoda: *Consistency Based Feature Selection*, 4th Pacific-Asia Conference, PADKK 2000, Kyoto, Japan, 2000