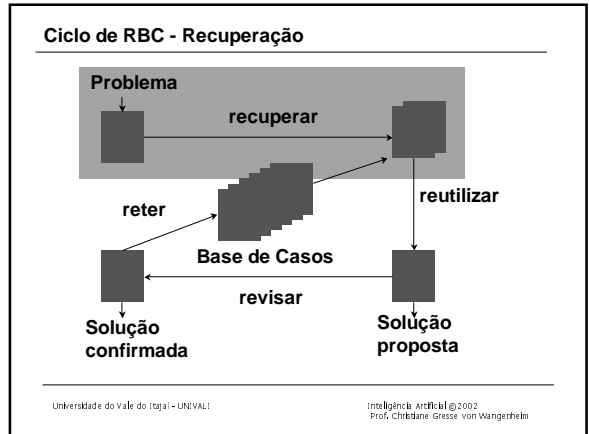


RBC-Recuperação

2002-2

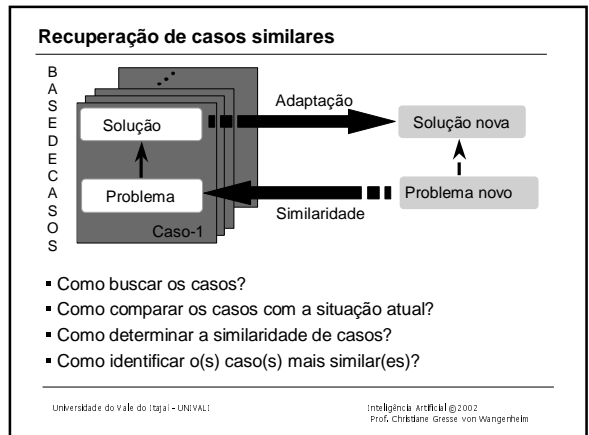
Christiane Gresse von Wangenheim
 Disciplina Inteligência Artificial
 UNIVALI



Como identificar casos uteis?

- Procura-se por caso(s) na base, que, na situação atual é **útil** para determinar a sua solução.
- O que significa um caso ser útil para solucionar um determinado problema?
- Hipótese: **Problemas similares possuem soluções semelhantes**
- ⇒ O critério *a posteriori* da **utilidade de soluções** passa a ser reduzido ao critério *a priori* **similaridade de descrições de problema**: Um caso é útil se ele é similar à questão atual.
- ⇒ Busca de casos **similares**
- Objetivo da similaridade:
 - Selecionar casos que possam ser facilmente adaptados para o problema atual
 - Selecionar casos que quase têm a mesma solução como o problema atual

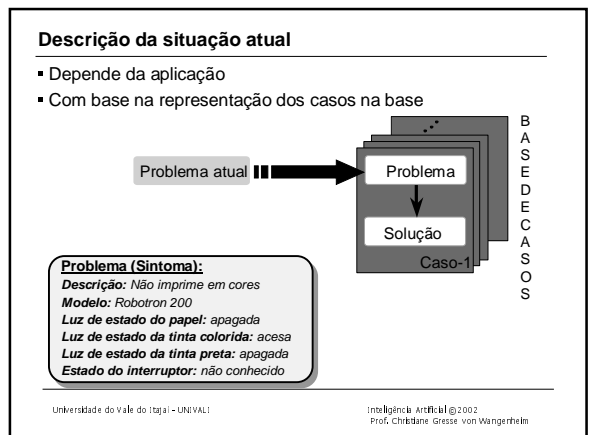
Universidade do Vale do Itajaí - UNIVALI
 Intelligência Artificial ©2002
 Prof. Christiane Gresse von Wangenheim



Tarefa central da recuperação

- **Dado:**
 - Base de casos $BC = \{C_1, \dots, C_n\}$ e uma medida de similaridade *sim*
 - Busca: *Q* (problema novo)
- **Queremos achar:**
 1. o caso mais similar C_i OU
 2. os m casos mais similares $\{C_1, \dots, C_m\}$ (ordenados ou sem ordem) OU
 3. todos os casos $\{C_1, \dots, C_m\}$ que têm pelo menos a similaridade sim_{min} com o problema novo Q
- **Processo de recuperação:**
 1. Descrição do problema/situação atual
 2. Busca na base de caso
 3. Comparação parcial dos casos da base com o problema atual
 4. Ordenação dos casos com base no valor da similaridade

Universidade do Vale do Itajaí - UNIVALI
 Intelligência Artificial ©2002
 Prof. Christiane Gresse von Wangenheim



Semântica da similaridade

- Grau da similaridade = utilidade/reusabilidade da solução
- Não existe uma similaridade absoluta - sempre depende da meta de recuperação na aplicação específica
 - dois carros são similares quando a velocidade máxima é similar?
 - dois carros são similares quando o preço é similar?
- A **meta de recuperação** explicitamente define o **objeto** a ser reutilizado, a **finalidade** de sua reutilização, a **tarefa** relacionada à reutilização, o **ponto de vista** específico e o **contexto** particular.
 - P.ex.: recupere o relatório de problema para o diagnóstico relativa ao conserto de impressoras do ponto de vista do pessoal do SAC na empresa IntelliPrinters.
- Meta da modelagem da similaridade: prover uma aproximação boa
 - perto da utilidade real
 - fácil de computar

Universidade do Vale do Itajaí - UNIVALI
 Intelligência Artificial @2002
 Prof. Christiane Gresse von Wangenheim

Modelar similaridade

- **Várias abordagens dependendo da representação de casos**
- **Medidas de similaridade:**
 - Funções para comparar dois casos *sim*: Caso x Caso → [0..1]
 Supõe (P₁, S₁) e (P₂, S₂) são dois casos e P o problema atual.
 Se $sim(P, P_1) \geq sim(P, P_2)$ então não preferimos a solução S₂ em cima da solução S₁ para o problema atual P.
 - Similaridades são geralmente normalizadas em uma faixa de 0 a 1, onde 0 é a dissimilaridade total e 1 a coincidência absoluta, ou através de porcentagens, onde 100% é um casamento exato.
- **Similaridade global vs. local**
 - Medidas de similaridade local: similaridade no nível de atributos
 - Medidas de similaridade global: similaridade no nível de casos
 - combinação de medidas de similaridade local
 - consideração de importância/pesos diferentes de atributos

Universidade do Vale do Itajaí - UNIVALI
 Intelligência Artificial @2002
 Prof. Christiane Gresse von Wangenheim

Similaridade global: Nearest Neighbor

- Ocorrências em uma base de casos são vistas como pontos em um espaço multidimensional.
- A distância espacial entre as respectivas representações dos caso reflete a similaridade entre estes.
- A busca reduz-se à determinação do vizinho geometricamente mais próximo, após definição de uma medida de distância d.

Universidade do Vale do Itajaí - UNIVALI
 Intelligência Artificial @2002
 Prof. Christiane Gresse von Wangenheim

Similaridade global: Nearest Neighbor

Caso 1:
 Modelo: Robotron Matrix 600
 Luz da tinta: vermelha ...

Caso 2:
 Modelo: Robotron 400
 Luz da tinta: verde ...

Situação atual Q:
 Modelo: Robotron 200
 Luz da tinta: vermelha ...

a distância de Q ao caso 1: $d_1 = X_1 + Y_1 = 3 + 0 = 3$
 a distância de Q ao caso 2: $d_2 = X_2 + Y_2 = 1 + 3 = 4$
 ⇒ caso 1 é o vizinho mais próximo.

Universidade do Vale do Itajaí - UNIVALI
 Intelligência Artificial @2002
 Prof. Christiane Gresse von Wangenheim

Medidas de similaridade global

- Exemplo: *nearest neighbor ponderada*
- Dado duas descrições de problemas C1, C2 com os atributos y₁, ..., y_p utilizado para representação

$$SIM(C1, C2) = \sum_{j=1}^p w_j sim_j(C1, C2)$$

- sim_j: Medida da similaridade local para atributo j
- peso w_j: indica a importância do atributo j para a determinação da similaridade
- peso w_j = 0 ⇔ atributo não considerado para recuperação
- normalização é realizada com a divisão do valor de similaridade pela soma total dos pesos dos índices

Universidade do Vale do Itajaí - UNIVALI
 Intelligência Artificial @2002
 Prof. Christiane Gresse von Wangenheim

Medidas de similaridade local

- Devem ser definidas em relação:
 - ao tipo específico de um atributo, e
 - no contexto de aplicação específico.
- Exemplos
 - Número
 - Símbolo binário
 - Símbolo (ordenado, não-ordenado, taxonômico)
 - String

Universidade do Vale do Itajaí - UNIVALI
 Intelligência Artificial @2002
 Prof. Christiane Gresse von Wangenheim

Medida de similaridade local: números

- Com base na diferença dos valores: $\text{sim}_f(x,y) = f(x-y)$
- Exemplo: $\text{sim}_f(x,y) = |x-y|$
- com f em geral:
 - $f: \mathbb{R} \rightarrow [0..1]$ oder $\mathbb{N} \rightarrow [0..1]$
 - $f(0) = 1$ (Reflexividade)
 - $f(x)$: decrescente monótono para $x > 0$ e crescente monótono para $x < 0$
- Exemplos:
 - Função escada: só quando um caso é completamente inútil antes de uma certo grau de similaridade
 - polinomial (fator 1: linear)
- Simetria/Assimetria
 - medida simétrica: (caso - busca)
 - medida assimétrica (caso-busca)
Exemplo: preço máximo

Universidade do Vale do Itajaí - UNIVALI

Inteligência Artificial @2002
Prof. Christiane Gresse von Wangenheim

Medida de similaridade local: símbolos ordenados

- Símbolos ordenados representam valores simbólicos em uma determinada ordem.
- Exemplo: «*febre*»:{*baixa*: temperatura entre 36C e 37C, *média*: temperatura entre 37C e 38.5C, *alta*: temperatura acima de 38.5C} em ordem crescente.
- Medida de similaridade local:
 - Assinar um valor Integer a cada símbolo preservando a ordem
- Exemplo:
 - baixa --> 1
 - média --> 2
 - alta --> 3
- sim_f : Uso das mesmas medidas do Tipo Número
- Normalmente, a distância entre estes valores ordinais será igual, mas pode-se também definir valores não equidistantes.

Universidade do Vale do Itajaí - UNIVALI

Inteligência Artificial @2002
Prof. Christiane Gresse von Wangenheim

Medida de similaridade local: símbolos não ordenados

- Símbolos não ordenados representam valores sem qualquer ordem definida.
- Exemplo: lista de destinos de uma agência de viagens: (*Rio de Janeiro, Brasília, Miami, Paris ...*)
- Tabela de similaridade: $\text{sim}_f(x,y) = s[x,y]$
- Tipo do símbolo $T_A = \{v_1, \dots, v_k\}$

s[x,y]	v ₁	v ₂	...	v _k
v ₁	s[1,1]	s[1,2]		s[1,k]
v ₂	s[2,1]	s[2,2]		s[2,k]
...				
v _k	s[k,1]	s[k,2]		s[k,k]

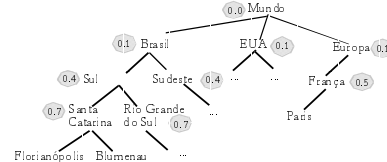
- Valores na diagonal = 1
- Medidas simétricas: Triângulo da matriz de cima = Triângulo da matriz de baixo

Universidade do Vale do Itajaí - UNIVALI

Inteligência Artificial @2002
Prof. Christiane Gresse von Wangenheim

Medida de similaridade local: taxonomias

- Uma taxonomia é uma árvore n-ária na qual os nodos representam valores simbólicos descrevendo o relacionamento entre os valores e sua posição na taxonomia.
- $\text{sim}_f(x,y)$
 - Assinar um valor de similaridade para cada nodo interno
 - Valores de similaridade crescente para os nodos sucessor
 - Similaridade entre dois nodos de folha é calculada pelo valor de similaridade do precedente comum mais próximo.



Universidade do Vale do Itajaí - UNIVALI

Inteligência Artificial @2002
Prof. Christiane Gresse von Wangenheim

Medida de similaridade local: strings

- Strings* descrevem valores de forma textual
- Exemplo «*problema da impressora*»: «*impressora não imprime em preto*».
- Calcular a similaridade entre *strings* considerando uma semântica razoável é uma tarefa extremamente difícil ⇨ substituir *strings* por valores simbólicos sempre que possível.
- correspondência exata**: dois *strings* são similares se são escritos da mesma forma
 - P.ex. $\text{sim}_f(\text{"printer"}, \text{"printer"}) = 1.0$; $\text{sim}_f(\text{"printer"}, \text{"print"}) = 0.0$.
- correção ortográfica**: compara o número de caracteres que são idênticos, ponderado pelo número total de caracteres no *string*-consulta.
 - P.ex. $\text{sim}_f(\text{"printer"}, \text{"print"}) = 5/7 = 0.7$
- contagem de palavras**: conta o número de palavras idênticas em dois casos, normalizado por meio da divisão pelo número total de palavras no *string*-consulta.
 - $\text{sim}_f(\text{"impressora não imprime preto"}, \text{"impressora não imprime texto azul"}) = 4/6 = 0.67$

Universidade do Vale do Itajaí - UNIVALI

Inteligência Artificial @2002
Prof. Christiane Gresse von Wangenheim

Processo de recuperação - 1

- O processo de recuperação de casos pode ser comparado com um problema de busca massiva.
- Idealmente, a medida de similaridade é aplicada a todos os casos gerado um conjunto-resposta **completo e correto**.
 - Completeza**: O método de recuperação é denominado completo, se toda relação de similaridade representada no modelo do sistema também se encontra no resultado deste método de recuperação.
 - Correção do método de recuperação**: Um método de recuperação de casos é correto, se uma relação de similaridade definida pelo método entre um caso e o problema atual também existe no conceito de similaridade desenvolvido para a aplicação.
- Mas, pelo problema de **eficiência** em grandes bases de casos: otimização de técnicas de busca otimizadas, que não analisam toda a base de casos para cada consulta, mas apenas um subconjunto desta considerado por alguma heurística

Universidade do Vale do Itajaí - UNIVALI

Inteligência Artificial @2002
Prof. Christiane Gresse von Wangenheim

Processo de recuperação - 2

- Várias abordagens:
 - Recuperação seqüencial
 - Recuperação de dois níveis
 - Recuperação usando árvores k-d
 - Recuperação usando redes
 - ...
- Abordagem depende do tamanho da base e da representação dos casos
- Organização da base de casos:
 - listas lineares (somente para bases pequenas)
 - estruturas indexadas para grandes bases:
 - *kd-trees*, redes de recuperação, etc.

Universidade do Vale do Itajaí - UNIVALI

Inteligência Artificial @2002
Prof. Christiane Gresse von Wangenheim

Recuperação seqüencial

```
TIPOS:
TipoCaso = ...
SimCaso = REGISTRO
    case: TipoCaso;
    similaridade: [0..1]
    FIM;
VARIAVEIS:
ListaCasoSim: VETOR [1..m] DE SimCaso
CaseBase: ARRAY [1..n] DE TipoCaso (* base de casos *)
Consulta: TipoCaso
```

Estrutura de dados

```
FUNÇÃO SelecRel (CaseBase, Consulta, m): ListaCasoSim Algoritmo
INICIO
ListaCasoSim[1..m].similaridade := 0
PARA i:=1 TO n FAÇA
    SE sim(Consulta, CaseBase[i]) > ListaCasoSim[m].similaridade
    ENTÃO insira CaseBase[i] em ListaCasoSim
RETORNE ListaCasoSim
FIM
```

Universidade do Vale do Itajaí - UNIVALI

Inteligência Artificial @2002
Prof. Christiane Gresse von Wangenheim

Características da recuperação seqüencial

- Complexidade: $O(n)$
- O processo é **completo e correto**, como o conceito de similaridade representado no sistema é aplicado de forma seqüencial a todos os exemplos de casos da base de casos:
- Desvantagens:
 - Lento, se a base é muito grande
 - Esforço da recuperação é independente da busca
 - Esforço da recuperação é independente do número dos casos a serem recuperados (m)
- Vantagens:
 - Implementação simples
 - Nenhuma estrutura de indexação necessária
 - Qualquer medida de similaridade pode ser utilizada

Universidade do Vale do Itajaí - UNIVALI

Inteligência Artificial @2002
Prof. Christiane Gresse von Wangenheim

Atividade curricular: Modelagem da recuperação

- Revisão da descrição do problema:
 - todos os atributos são relevantes para a recuperação de casos similares?
 - A descrição inclui todos os atributos relevantes?
- Definição das medidas de similaridade local para cada tipo de um atributo utilizado na recuperação
- Definição da medida de similaridade global (inclusive definição do peso para cada atributo)

Universidade do Vale do Itajaí - UNIVALI

Inteligência Artificial @2002
Prof. Christiane Gresse von Wangenheim