

# Análise Forense de Documentos Digitais

*Prof. Dr. Anderson Rocha*

[anderson.rocha@ic.unicamp.br](mailto:anderson.rocha@ic.unicamp.br)

<http://www.ic.unicamp.br/~rocha>

---

Reasoning for Complex Data (RECOD) Lab.  
Institute of Computing, Unicamp

Av. Albert Einstein, 1251 – Cidade Universitária  
CEP 13083-970 • Campinas/SP – Brasil

---

# Organização

# Organização

- ▶ Conceitos de Imagem Digital
- ▶ Operações com Imagens
- ▶ Aprendizado de Máquina

# Organização

- ▶ Aprendizado de Máquina
  - Supervisionado
  - Não-Supervisionado
  - Semi-Supervisionado
- ▶ Avaliação e Comparação de Métodos

**Imagem**

# Imagem

- ▶ De acordo com [Gomes & Velho 1996], para trabalharmos com imagens, devemos estabelecer um **universo matemático** no qual seja possível definir diversos modelos abstratos destas
- ▶ Em seguida, precisamos criar um **universo de representação** onde procuramos esquemas que permitam uma representação discreta desses modelos

# Imagem

- ▶ O objetivo da representação discreta desses modelos é **codificar a imagem no computador**
- ▶ Quando observamos uma fotografia, ou uma cena no mundo real, recebemos de cada ponto do espaço um impulso luminoso que associa uma informação de cor a esse ponto

# Imagem

- ▶ Nesse sentido, podemos definir uma imagem contínua (não discreta) como a aplicação

$$\mathcal{I} : \mathcal{U} \rightarrow \mathcal{C}$$

onde  $\mathcal{U} \subset \mathbb{R}^3$  é uma superfície e  $\mathcal{C}$  é um espaço vetorial

- ▶ Na maioria das aplicações,  $\mathcal{U}$  é um **subconjunto plano** e  $\mathcal{C}$  é um **espaço de cor**



# Imagem

- ▶ A função  $\mathcal{I}$  na definição é chamada de **função imagem**
- ▶ O conjunto  $\mathcal{U}$  é chamado **suporte da imagem**
- ▶ O conjunto de valores de  $\mathcal{I}$ , que é um subconjunto de  $\mathcal{C}$ , é chamado de conjunto de **valores da imagem**

# Imagem

- ▶ Quando  $\mathcal{C}$  é um espaço de cor de dimensão 1, dizemos que a imagem é monocromática ou em tons de cinza
- ▶ A representação mais comum de uma imagem espacial consiste em tomar um subconjunto discreto  $\mathcal{U}' \subset \mathcal{U}$  do domínio da imagem, um espaço de cor  $\mathcal{C}$  associado a um dispositivo gráfico e representar a imagem pela amostragem da função imagem  $\mathcal{I} \rightarrow \mathcal{U}'$

# Imagem

- ▶ Cada ponto  $(x_i, y_i)$  do subconjunto discreto  $\mathcal{U}'$  é chamado de elemento da imagem ou **pixel**
- ▶ Para a representação em computador, devemos também trabalhar com modelos onde a função imagem  $\mathcal{I}$  toma valores em um subconjunto discreto do espaço de cor  $\mathcal{C}$
- ▶ Esse processo de discretização é chamado de **quantização**

# Imagem

- ▶ O caso mais utilizado de discretização espacial de uma imagem consiste em tomar o domínio como sendo um **retângulo** e discretizar esse retângulo usando os pontos de um **reticulado bidimensional**
- ▶ Dessa forma a imagem pode ser representada de forma matricial por uma matriz

$$A^{(m \times n)} = (a_{ij} = (\mathcal{I}(x_i, y_j)))$$

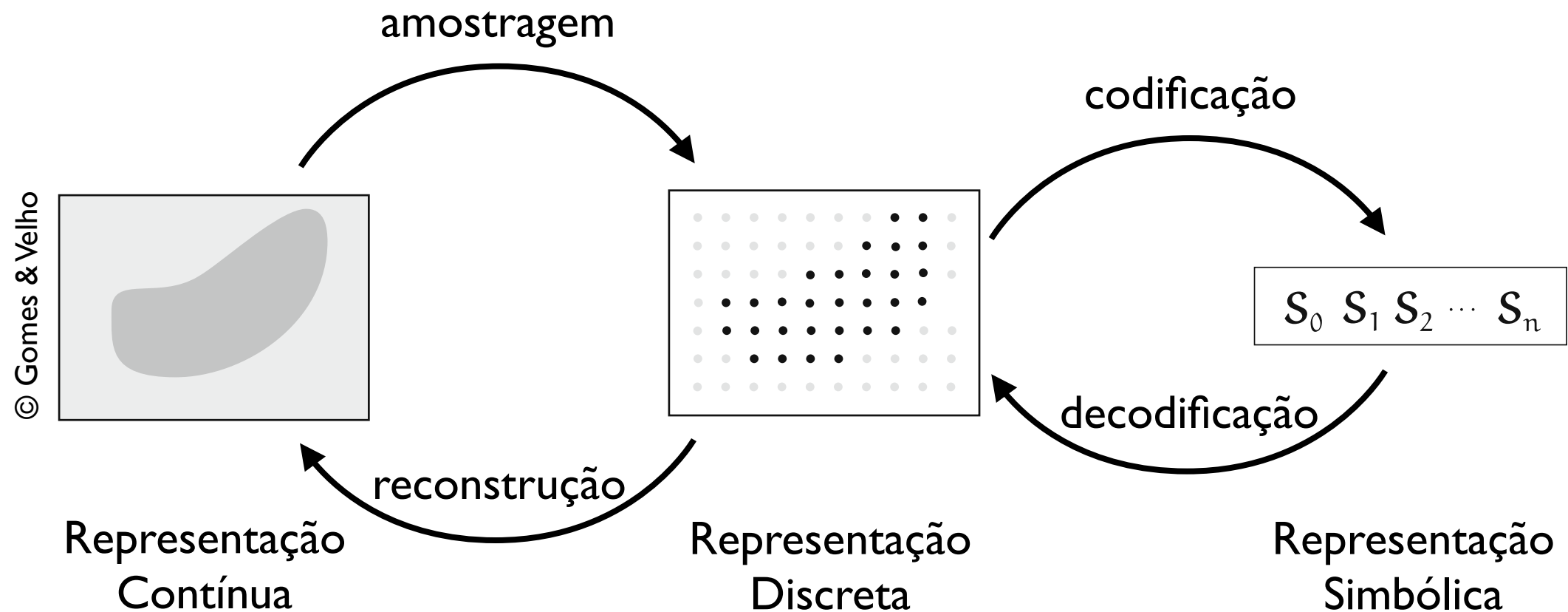
# Imagem

- ▶ Cada elemento  $a_{ij}, i = 1, \dots, m$  e  $j = 1, \dots, n$  da matriz representa o valor da função imagem  $\mathcal{I}$  no ponto de coordenadas  $(x_i, y_j)$  do reticulado
- ▶ Dessa forma, cada ponto  $a_{ij}$  é um vetor do espaço de cor representando a cor do *pixel* na coordenada  $(i, j)$  da imagem

# Imagem

- ▶ Se cada ponto possui três valores associados e cada valor precisa de oito *bits* para ser representado, então cada *pixel* dessa imagem pode ser representado com 24 *bits*
- ▶ A imagem é dita de 24 *bits*
- ▶ Se cada *pixel* também codifica **transparência**, a imagem tem um quarto canal, chamado alfa, tornando-se uma imagem de 32 *bits*

# Imagem



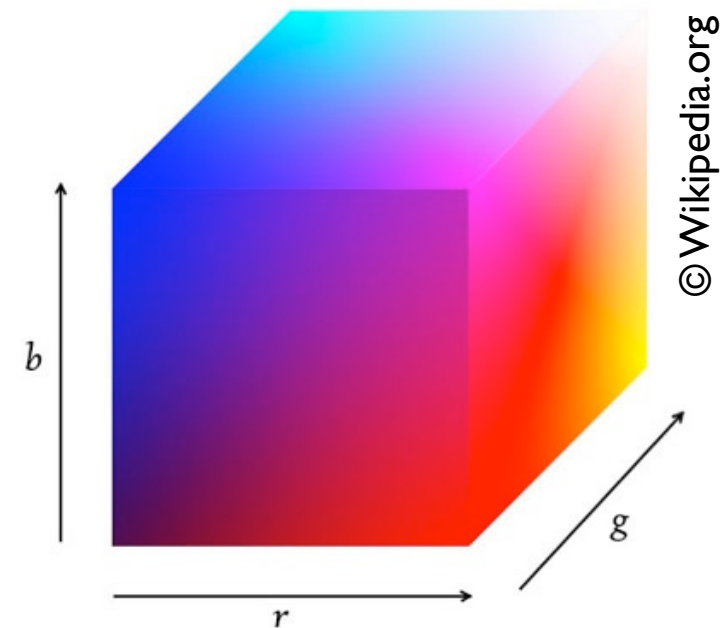
# Espaços de Cor

- ▶ O espaço de cor pode variar de acordo com o dispositivo de exibição (e.g., monitor, impressora)
- ▶ Espaços de cor
  - RGB (Vermelho, Verde, Azul)
  - CMYK (Ciano, Magenta, Amarelo, Preto)
  - HSV (Matiz, Saturação e Brilho)
  - etc.

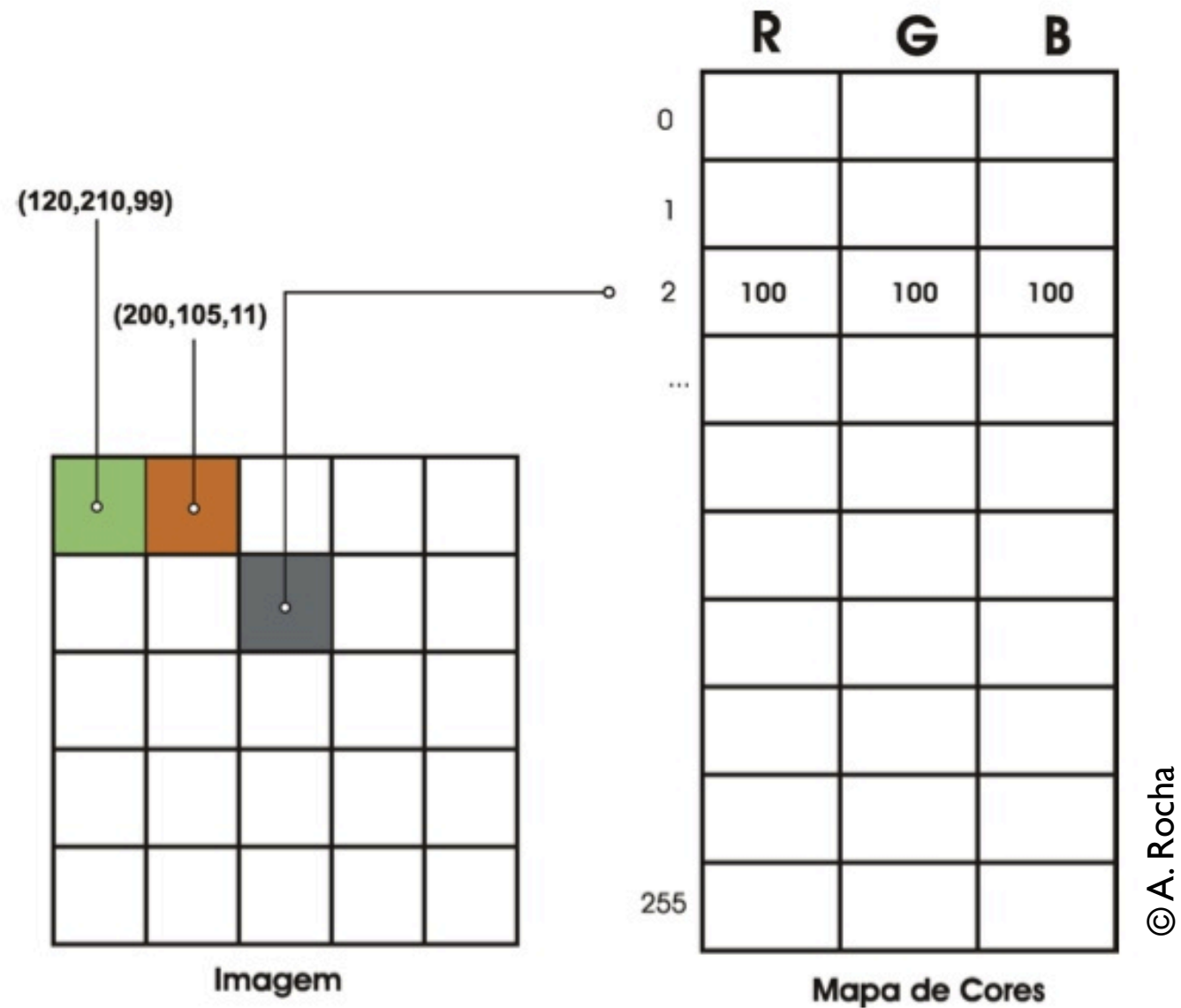


# Espaço de cor RGB

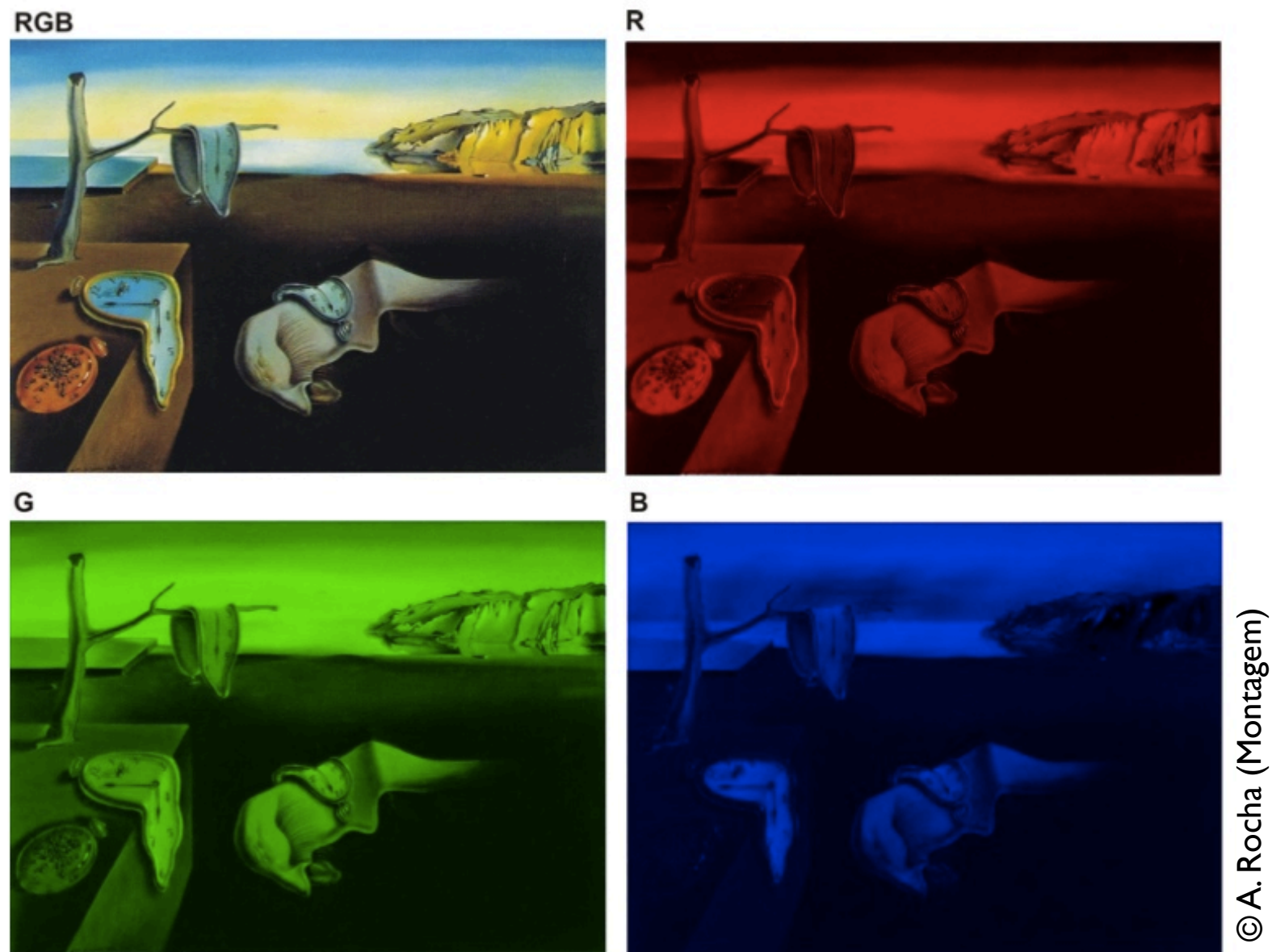
- ▶ O propósito principal do sistema RGB é a reprodução de cores em dispositivos eletrônicos
  - monitores de TV e computador
  - *datashows*
  - *scanners*
  - câmeras digitais
  - fotografia tradicional



# Imagem



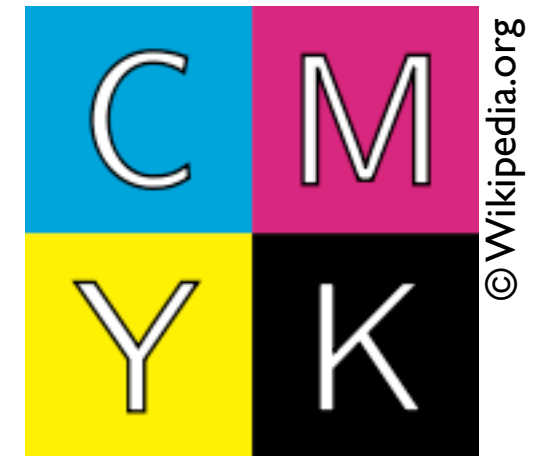
# Imagem



\*The Persistence of Memory by Salvador Dali

# Espaço de cor CMYK

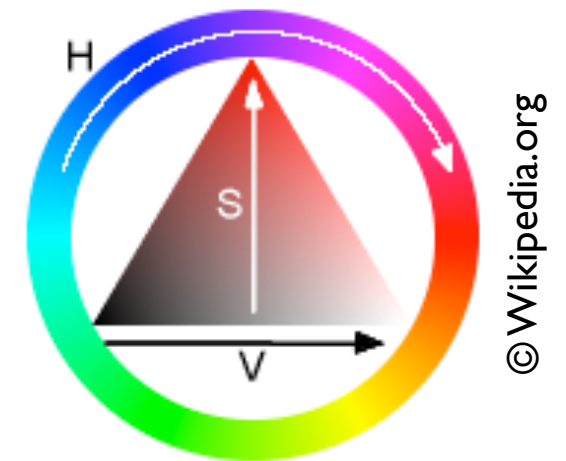
- ▶ Modelo de cores subtrativas
- ▶ Contraposição ao RGB
- ▶ Adequado para impressoras
- ▶ *K* vem de *keyed* (alinhamento) da placa de impressão de cor preta com as outras



© Wikipedia.org

# Espaço de cor HSV

- ▶ **Matiz** (tonalidade): verifica o tipo da cor (abrange todas as cores do espectro)
- ▶ **Saturação** (pureza): valores baixos são próximos do cinza. Valores altos são próximos da cor pura
- ▶ **Brilho**: define o brilho (intensidade) da cor



© Wikipedia.org

# **Operações com Imagens**

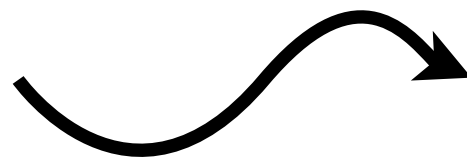
# Quantização

- ▶ Mapeamento dos números reais em valores discretos
- ▶ Tipicamente utiliza-se *bytes* (256 valores) ou inteiros curtos (65536 valores)



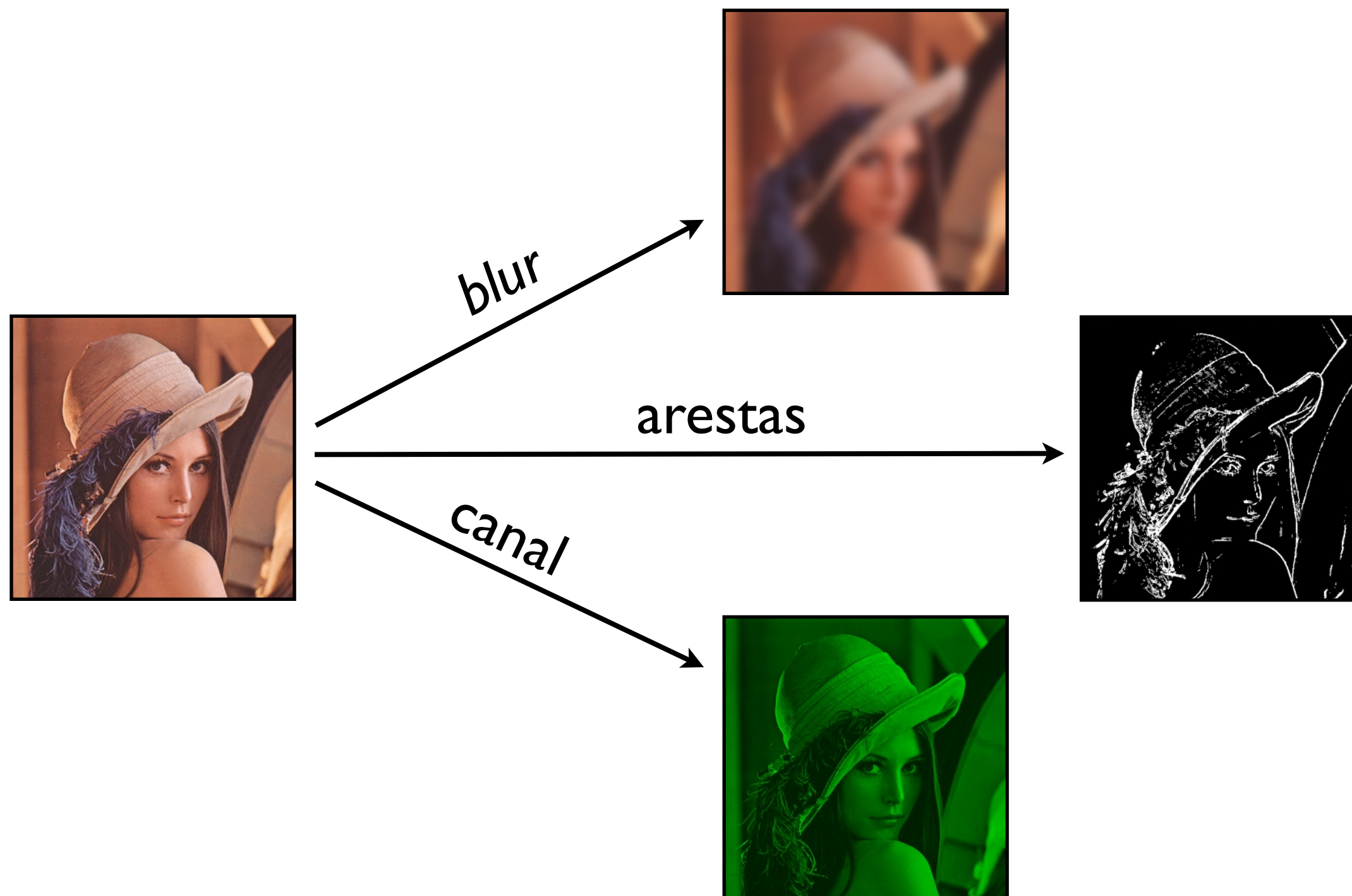
# Warping

- Modifica o “domínio” da função de imagem.





# Transformações - Atributos

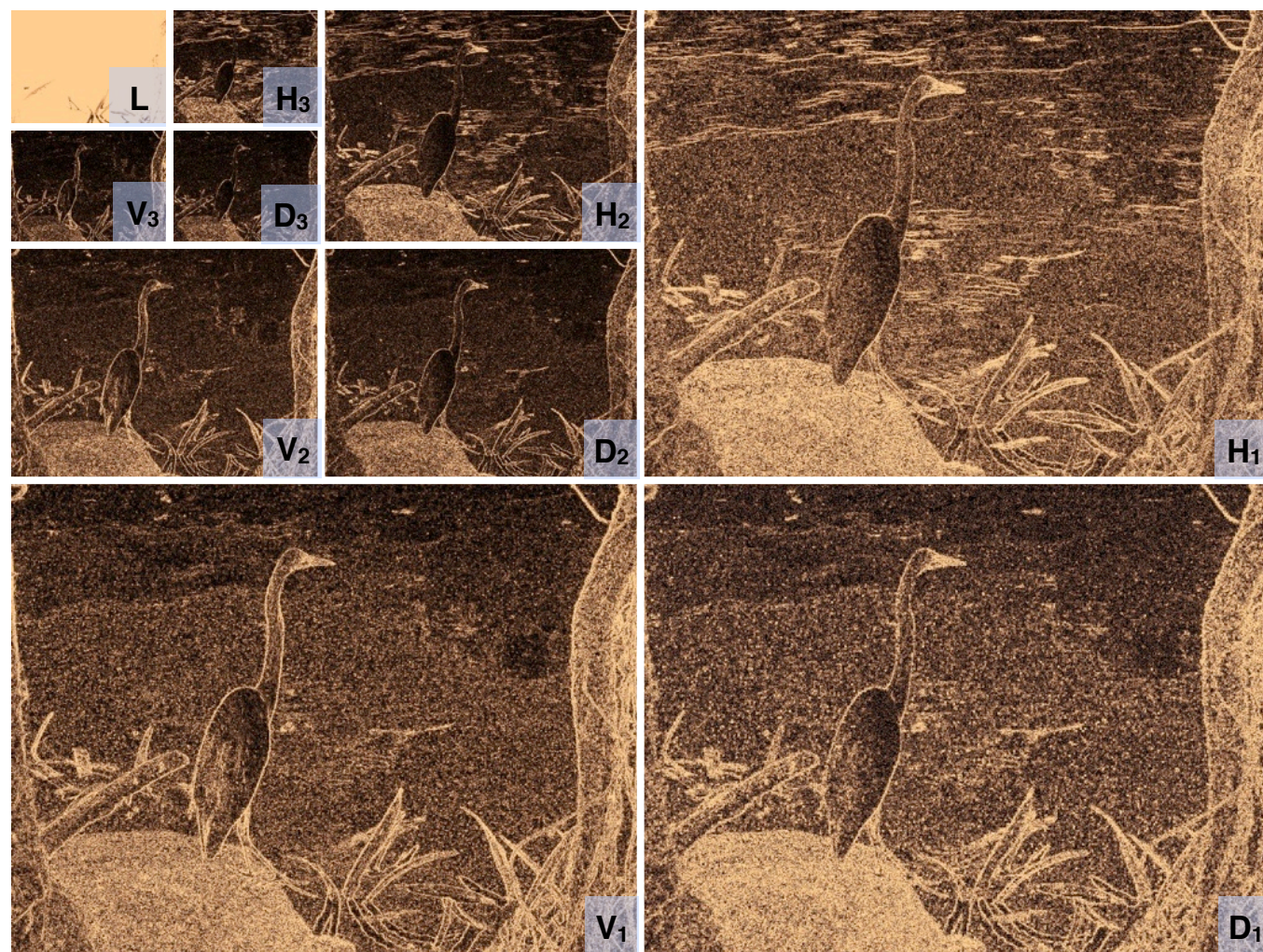


# Decomposição em Canais de Cores

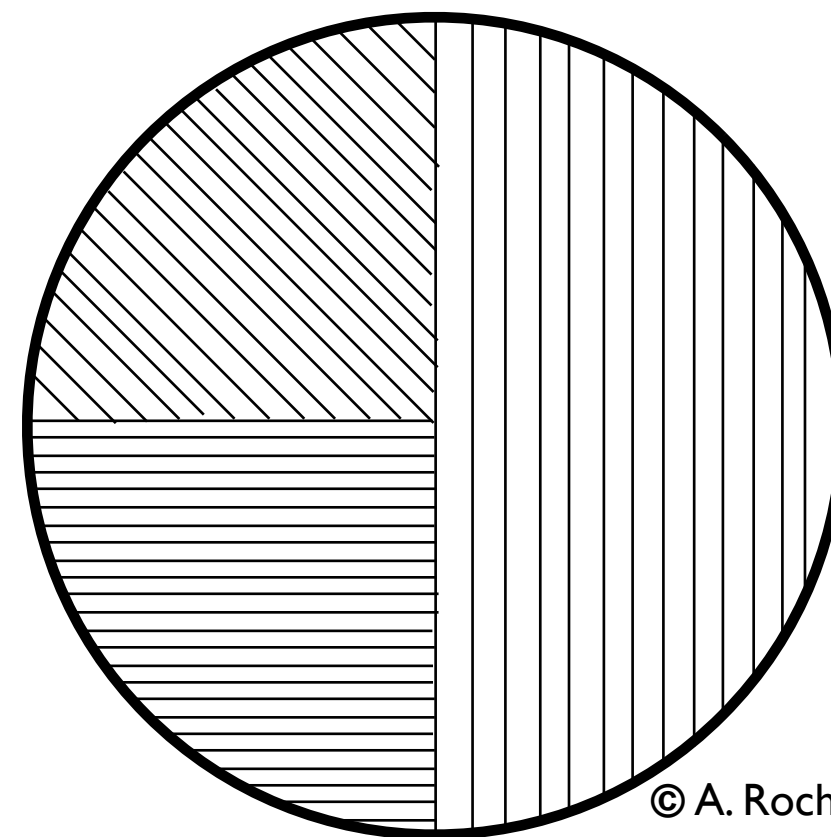
- ▶ Quando separamos a imagem em suas cores básicas representadas no espaço de cores  
 $\mathcal{C}' \in \mathcal{C}$
- ▶ Se o espaço de cores utilizado é um espaço RGB, temos os componentes vermelho (*Red*), verde (*Green*), e azul (*Blue*);



# Decomposição Wavelet



© A. Rocha

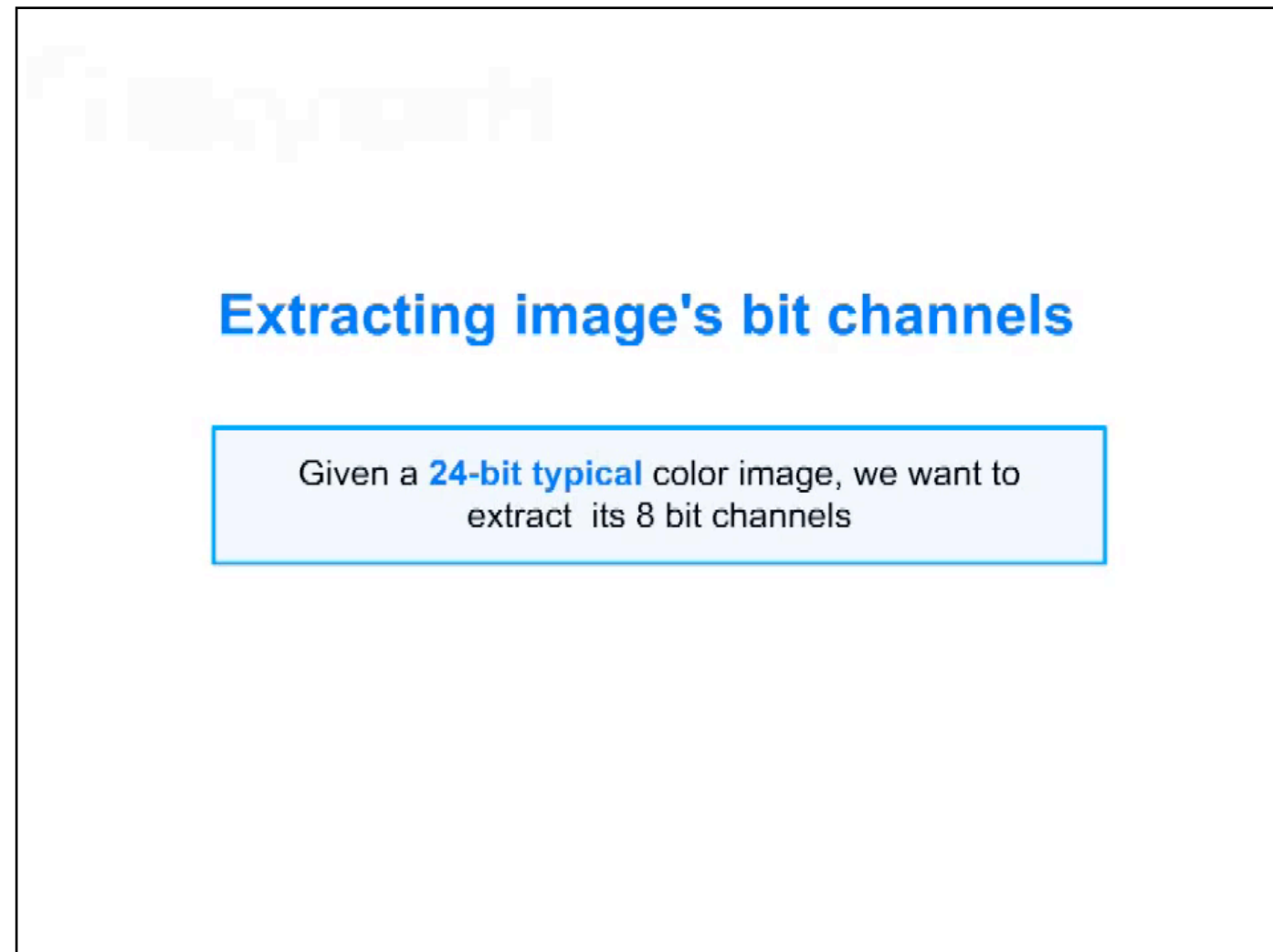


© A. Rocha

# Decomposição em Planos de Bits

- ▶ Quando decompomos a imagem em seus planos de *bits*
- ▶ Por exemplo, após a decomposição da imagem de 24 *bits* em seus três canais de cores (R,G,B), podemos ainda, fazer uma decomposição por planos de *bits*.
- ▶ Cada canal de cor possui 8 *bits* e possui 8 planos de *bits* por canal de cor

# Decomposição em Planos de Bits



\* Decomposição da imagem em canais de *bits*

# Nomenclaturas



# Nomenclatura

- ▶ Diferentes áreas tem nomes distintos para coisas parecidas
  - Aprendizado de Máquina
  - Reconhecimento de Padrões
  - Aprendizado Estatístico
  - Mineração de Dados

# **Aprendizado de Máquina**



# Aprendizado de Máquina

- ▶ Aprendizado de Máquina é uma área da Inteligência Artificial concentrada no desenvolvimento de técnicas que permitem que computadores sejam capazes de aprender com a experiência [Mitchell 1997]
- ▶ Extração de informações e extrapolação do conhecimento a partir de dados

# Aprendizado de Máquina

- ▶ Alguns problemas que utilizam aprendizado de Máquina [Mitchell 1997] [Friedman et al. 2001]
  - reconhecimento de caracteres
  - reconhecimento da fala
  - predição de ataques cardíacos
  - detecção de fraudes em cartões de créditos

# Aprendizado de Máquina

- ▶ Na solução desses problemas, podemos ter classificadores fixos ou baseados em aprendizado, que, por sua vez, pode ser supervisionado ou não-supervisionado [Friedman et al. 2001]

# Definição – Classificadores

- ▶ Podemos ver um classificador, matematicamente, como um mapeamento a partir de um espaço de características  $X$  para um conjunto discreto de rótulos (*labels*)  $Y$
- ▶ Em IA, um classificador de padrões é um tipo de motor de inferência que implementa estratégias eficientes para computar relações de classificação entre pares de conceitos ou para computar relações entre um conceito e um conjunto de instâncias  
[Duda et al. 2000]

# Classificadores

- ▶ Classificadores podem ser
  - Supervisionados
  - Semi-Supervisionados
  - Não-Supervisionados

# Classificadores

- ▶ Classificadores supervisionados consistem em técnicas em que procuramos estimar uma função de classificação  $f$  a partir de um conjunto de treinamento
- ▶ O conjunto de treinamento consiste de pares de valores de entrada  $X$ , e sua saída desejada  $Y$   
[Friedman et al. 2001]

# Classificadores

- ▶ Valores observados no conjunto  $X$  são denotados por  $x_i$ , isto é,  $x_i$  é a  $i$ -ésima observação em  $X$
- ▶ O número de variáveis que constituem cada uma das entradas em  $X$  é  $p$
- ▶ Assim,  $X$  tem  $n$  observações, chamados de vetores de características

# Classificadores

- ▶ Cada vetor de entrada é composto por  $p$  graus de liberdade (dimensões ou variáveis)
- ▶ A saída da função  $f$  pode ser um valor contínuo (regressão) ou pode predizer a etiqueta (*label*) de um objeto de entrada (*classificação*)



# Classificadores

- ▶ A tarefa do aprendizado é prever o valor da função para qualquer objeto de entrada que seja válido após ter sido suficientemente treinado com um conjunto de exemplos [Bishop 2006]
- ▶ Alguns exemplos de classificadores supervisionados são
  - *Support Vector Machines*
  - *Linear Discriminant Analysis,*
  - *Boosting*

# Aprendizado Não-Supervisionado

- ▶ Um outro grupo de técnicas de aprendizado, não utilizam exemplos de treinamento marcados (classe conhecida)
- ▶ Conhecidos como técnicas para aprendizado não-supervisionado
- ▶ Esta forma de aprendizado, na maioria das vezes, trata o seu conjunto de entrada como um conjunto de variáveis aleatórias

# Aprendizado Não-Supervisionado

- ▶ Um modelo de distribuição conjunta (*joint distribution model*) é então construído para a representação dos dados
- ▶ Desta forma, o objetivo deste aprendizado é avaliar como os dados estão organizados e agrupados [Friedman et al. 2001]
- ▶ Técnicas de Maximização de Esperança [Baeza-Yates 2003], por exemplo, podem ser utilizadas para aprendizado não-supervisionado

# Aprendizado Semi-Supervisionado

- ▶ Um outro grupo de técnicas de aprendizado envolve abordagens mistas
  - Supervisionado
  - Não Supervisionado
- ▶ São as técnicas **Semi-Supervisionadas**

# **Modelagem de Problemas**

# Modelagem de Problemas

- ▶ Problemas são descritos por variáveis
- ▶ Dois tipos
  - Reais
  - Categóricas

# Modelagem de Problemas

- ▶ Como transitar entre os dois tipos de variáveis?
- ▶ É possível “converter” uma representação em outra?

# Modelagem de Problemas

- ▶ **Simplicidade vs. Complexidade**
- ▶ O que é realmente importante?
- ▶ Precisamos realmente de todos os dados possíveis para tomar uma decisão?



# Modelagem de Problemas

- ▶ Dimensão do vetor de características tem efeitos colaterais importantes:
- ▶ Dimensão alta
  - Distâncias médias ficam grandes
  - Dados ficam esparsos
- ▶ **Maldição da Dimensionalidade**

# **Aprendizado Supervisionado**

*(Primeiros Passos)*

# Aprendizado Supervisionado

- ▶ Dados para Aprendizado Supervisionado
- ▶ “Give me more data”
- ▶ Classificação vs. Regressão

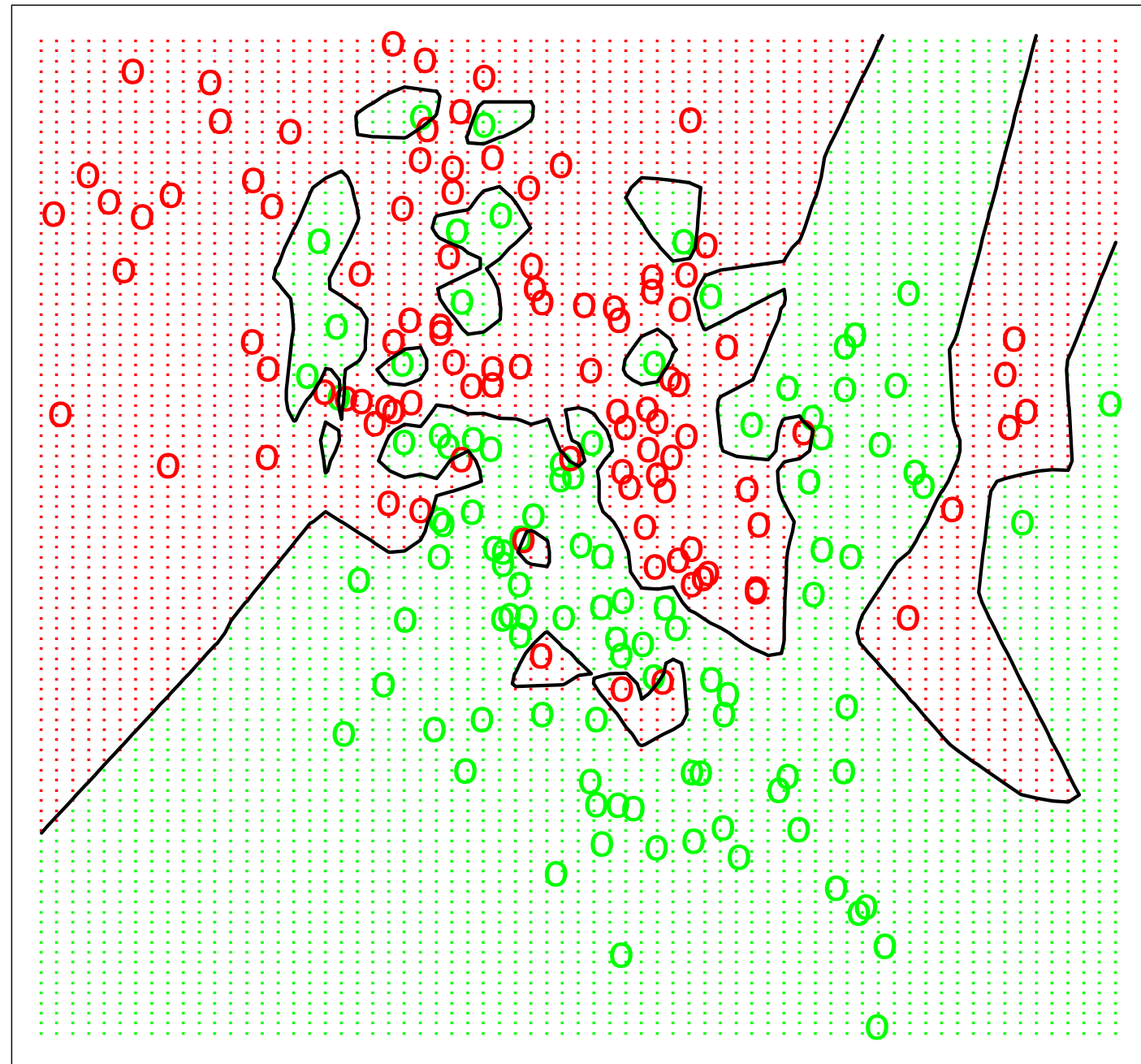
# Aprendizado Supervisionado

- ▶ Será que quanto mais complexo nosso modelo de “predição” melhor o resultado?

# Exemplo – KNN

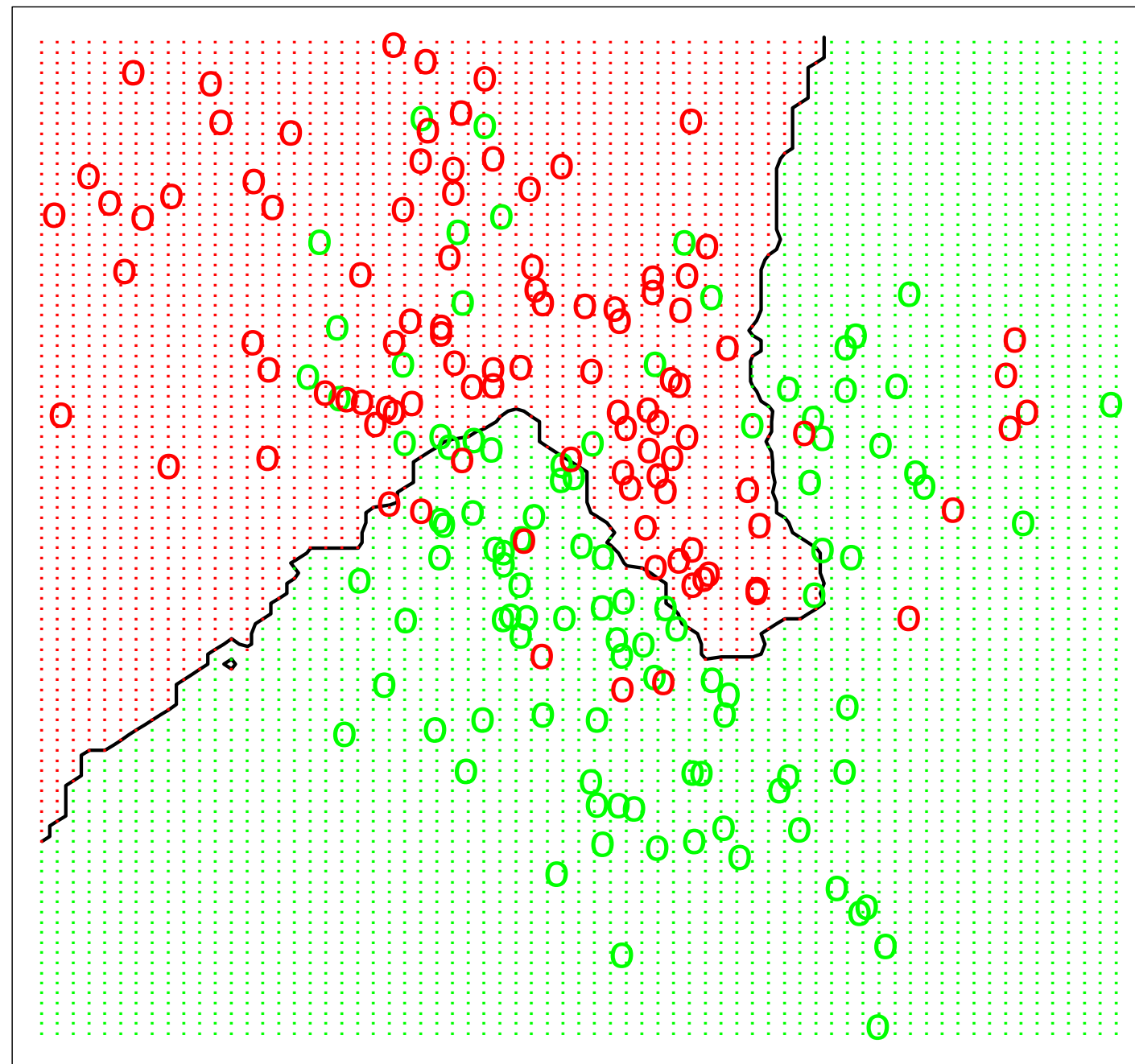
- ▶ K-Vizinhos mais Próximos (KNN)
- ▶ Um exemplo de técnica baseada em instâncias.
- ▶ Não há “aprendizado”
  - decisões são feitas para cada instância

# KNN – $k = 1$



© J. Friedman et al.

# KNN – $k = 15$



© J. Friedman et al.

# **Avaliação e Comparação**



# Avaliação e Comparação

- ▶ Viés e Variância
- ▶ Treinamento e Teste
- ▶ Matriz de Confusão
- ▶ Métricas e Critérios

# Avaliação e Comparação

- ▶ Conjuntos de validação e teste
- ▶ Validação cruzada

# Curvas ROC

## ► Especificidade

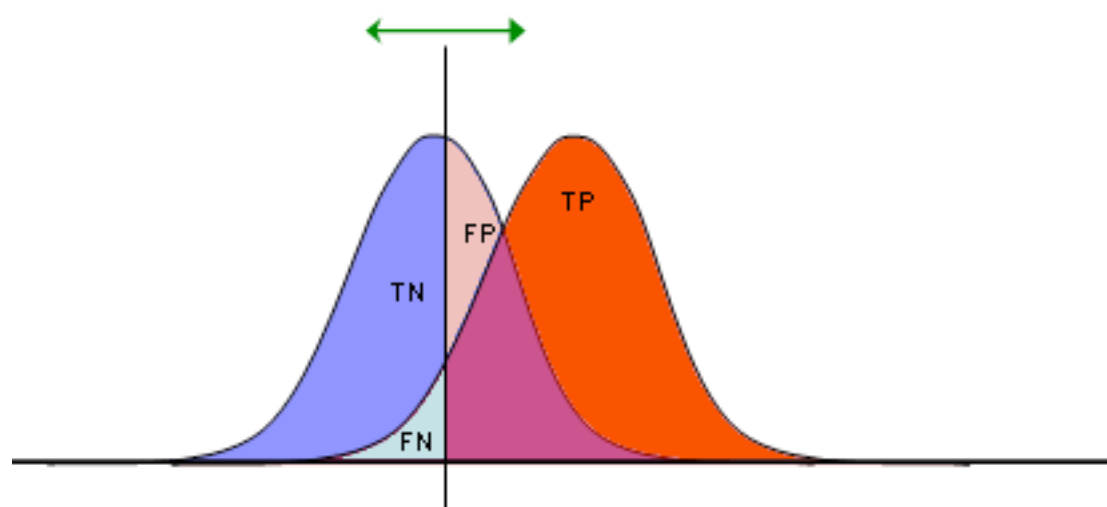
- $E = TN / (TN + FP)$

## ► Sensitividade

- $S = TP / (TP + FN)$

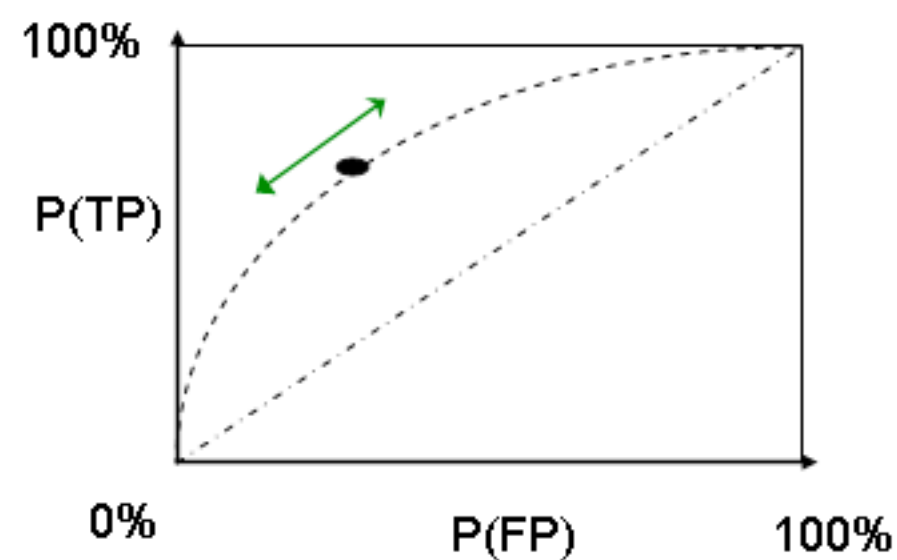
## ► (Sensitividade) vs. (1 - Especificidade) = Curva Característica de Operação (ROC)

# Curvas ROC



TP	FP
FN	TN
1	1

© Wikipedia.org



# Referências

# Referências

1. [Baeza-Yates 2003] **R. Baeza-Yates**. *Clustering and Information Retrieval*. Kluwer Academic Publishers. 1 edition.
2. [Bishop, 2006] **C. M. Bishop**. *Pattern Recognition and Machine Learning*. Springer, 1 edition, 2006.
3. [Duda et al. 2001] **R. O. Duda, P. E. HART and D. G. STORK**. *Pattern Classification*. Wiley-Interscience, 2, 2000.
4. [Friedman et al. 2001] **J. Friedman, T. Hastie, and R. Tibshirani**. *The Elements of Statistical Learning*. Springer, 1 edition, 2001.
5. [Gomes & Velho, 1996] **J. Gomes L. Velho**. *Computação Gráfica: Imagem*. IMPA-SBM, 1.
6. [Gonzalez & Woods, 2007] **R. Gonzalez and R. Woods**. *Digital Image Processing*. Prentice-Hall, 3 edition.
7. [Mitchell 1997] **T. M. Mitchell**. *Machine Learning*. McGraw-Hill, 1 edition, 1997.

---

***Obrigado!***

---