

UNIVERSIDADE FEDERAL DO MARANHÃO
CENTRO DE CIÊNCIAS EXATAS E TECNOLOGIA-CCET
CURSO DE CIÊNCIA DA COMPUTAÇÃO

JEFFERSON ALVES DE SOUSA

**CLASSIFICAÇÃO DE IMAGENS DE MASSAS EM MAMOGRAFIA USANDO LBP,
ÍNDICE DE DIVERSIDADE E SVM**

São Luís

2015

JEFFERSON ALVES DE SOUSA

**CLASSIFICAÇÃO DE IMAGENS DE MASSAS EM MAMOGRAFIA USANDO LBP,
ÍNDICE DE DIVERSIDADE E SVM**

Monografia apresentada ao curso de Ciência da Computação da Universidade Federal do Maranhão como parte dos requisitos necessários para obtenção do grau de Bacharel em Ciência da Computação.

Orientador: Prof. Dr. Anselmo Cardoso de Paiva

São Luís
2015

Sousa, Jefferson Alves de.

Classificação de imagens de massas em mamografia usando LBP, índice de diversidade e SVM/ Jefferson Alves de Sousa. – São Luís, 2015.

47 f.

Impresso por computador (fotocópia).

Orientador: Anselmo Cardoso de Paiva.

Monografia (Graduação) – Universidade Federal do Maranhão, Curso de Ciência da Computação, 2015.

1. Mamografia. 2. Local Binary Pattern. 3. Índices de Diversidade de Gleason e Menhinick. I. Título.

CDU 004.932:618.19

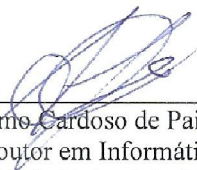
JEFFERSON ALVES DE SOUSA

**CLASSIFICAÇÃO DE IMAGENS DE MASSAS EM MAMOGRAFIA USANDO LBP,
ÍNDICE DE DIVERSIDADE E SVM**

Monografia apresentada ao curso de Ciência da Computação da Universidade Federal do Maranhão como parte dos requisitos necessários para obtenção do grau de Bacharel em Ciência da Computação.

Aprovada em: 13 / 01 / 2015

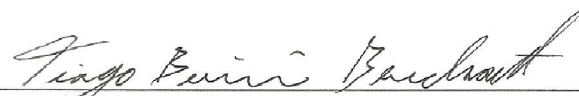
BANCA EXAMINADORA



Prof. Dr. Anselmo Cardoso de Paiva (Orientador)
Doutor em Informática
Universidade Federal do Maranhão



Prof. Dr. Geraldo Braz Júnior
Doutor em Engenharia de Eletricidade
Universidade Federal do Maranhão



Prof. Dr. Tiago Bonini Borchart
Doutor em Computação
Universidade Federal do Maranhão

Aos meus pais, pela confiança, apoio e incentivo.
Às minhas irmãs, Erika e Jessika, pela ajuda e
compreensão nos momentos difíceis.

AGRADECIMENTOS

Deus em primeiro lugar.

À minha família, em especial a meu pai, minha mãe e minhas irmãs, que acreditaram no meu potencial, e sempre estiveram dispostos a me ajudar em tudo que fosse necessário.

Aos professores do curso de Ciência da Computação da Universidade Federal do Maranhão, pelo conhecimento que me proporcionaram.

Ao professor Anselmo, meu orientador, pela paciência, dedicação e conhecimentos passados.

À professora Simara e ao professor Aristófanés pelas dicas e orientação em projetos durante a graduação.

Aos amigos do curso, principalmente Caio Eduardo, Giovanni, João e Whesley, por esse período de convivência, troca de conhecimento e experiências tanto boas quanto más, que contribuíram para o meu desenvolvimento.

À minha prima Queli, pela ajuda na finalização deste trabalho.

Por fim, agradeço a todos que direta ou indiretamente contribuíram para a realização deste trabalho.

“O sucesso nasce do querer, da determinação e persistência em se chegar a um objetivo. Mesmo não atingindo o alvo, quem busca e vence obstáculos, no mínimo fará coisas admiráveis”
(José de Alencar)

LISTA DE FIGURAS

Figura 1. Etapas do processamento de imagens. Fonte: (GONZALES, 1992).....	17
Figura 2. Exemplo de uma co-ocorrência dos níveis de cinza i e j , com vizinhança $d = 4$, alinhados na horizontal ($\theta = 0$). Fonte: (CARVALHO, 2012)	20
Figura 3. (a) Imagem de $M \times N$ pixels. (b) Matriz de Co-Ocorrência de Níveis de Cinza da imagem ($d = 2, \theta = 0$). Fonte: (CARVALHO, 2012).....	21
Figura 4. Exemplo de uma corrida de nível de cinza i , de comprimento 10 e direção horizontal. Fonte: (CARVALHO, 2012)	21
Figura 5. (a) Imagem de $M \times N$ pixels. (b) Matriz de Comprimento de Corrida de Níveis de Cinza da imagem ($\theta = 0$). Fonte: (CARVALHO, 2012)	22
Figura 6. Lacuna de nível de cinza g , de comprimento l e direção horizontal. Fonte: (CARVALHO, 2012).....	23
Figura 7. (a) Imagem de $M \times N$ pixels. (b) Matriz de Comprimento de Lacuna de Níveis de Cinza da imagem ($\theta = 0$). Fonte: (CARVALHO, 2012)	24
Figura 8. Exemplo do calculo do LBP. Fonte: (ROCHA, 2012).....	25
Figura 9. Separação entre duas classes através de hiperplanos. Fonte: (NASCIMENTO, 2012).	28
Figura 10. Vetores de Suporte (destacado por círculos). Fonte: (NASCIMENTO, 2012).....	29
Figura 11. Etapas da metodologia proposta.	33
Figura 12. (a) Mamografia com nódulo benigno selecionado. (b) Mamografia com nódulo maligno selecionado. Fonte: (BRAZ, 2008).	34
Figura 13. (a) Imagem original e seu histograma. (b) Imagem realçada e seu histograma. (c) Imagem suavizada e seu histograma. Fonte: (ROCHA, 2014).	35
Figura 14. Cálculo da matriz GLCM para $\theta = 0^\circ$ e $d = 2$. (a) ROI 5 x 5. (b) Ocorrência de pares de LBPs de mesmo valor. (c) Ocorrência de pares de LBPs de valores diferentes. Fonte: (ROCHA, 2014).....	36
Figura 15. Cálculo da matriz GLRLM para $\theta = 0^\circ$. (a) ROI 5 x 5. (b) Ocorrência de corridas de LBPs de comprimento $k = 3$. Fonte: (ROCHA, 2014).....	37
Figura 16. Cálculo da matriz GLGLM para $\theta = 0^\circ$. (a) ROI 5 x 5. (b) Ocorrência de lacunas de LBPs de comprimento $k = 2$. Fonte: (ROCHA, 2014).....	37
Figura 17. Fluxo de atividade da etapa de classificação. Fonte: (BRAZ, 2008).....	38

LISTA DE TABELAS

Tabela 1: Resultados do Índice de Gleason utilizando SVM	41
Tabela 2: Resultados do Índice de Gleason utilizando Adaboost.M1.....	41
Tabela 3: Resultados do Índice de Menhinick utilizando SVM	42
Tabela 4: Resultados do Índice de Menhinick utilizando Adaboost.M1.....	42
Tabela 5: Comparação com alguns trabalhos referentes à classificação de massas em imagens de mamografias em maligno e benigno.	43

RESUMO

Este trabalho tem o objetivo de investigar a aplicação de técnicas de análise de textura e reconhecimento de padrões para diagnóstico de câncer de mama, cujo objetivo é dar ao especialista um maior suporte ao diagnóstico do câncer de mama. Busca-se utilizar somente a textura para caracterizar o padrão maligno e benigno ao invés das características do contorno das massas, já que tais características nem sempre são nítidas nas imagens, pois pode existir desde a sobreposição de achados como, por exemplo, massas e calcificações até lesões que não têm contorno bem definido, impedindo a visualização das mesmas e contribuindo para um número maior de biopsias com resultados negativos.

Assim, este trabalho se propõe a estudar técnicas de análise de textura, tais como: *Local Binary Pattern* e Índices de Diversidade de *Gleason e Menhinick*, pois acredita-se que tais técnicas possam produzir boas características de textura que discriminem as regiões de massas nas mamografias digitalizadas entre malignas e benignas, visto que o sucesso da etapa de classificação depende muito das características geradas. As características produzidas serão submetidas como entrada para o processo de classificação supervisionada usando SVM e a estratégias de combinações de modelos *Ensemble*. Em ambos os casos uma parte das amostras será usada para a etapa treinamento do classificador. Esta etapa cria um padrão sobre as medidas extraídas. A outra parte, totalmente desconhecida da etapa de treinamento, é utilizada para fazer os testes e a validação dos resultados. Por último segue a etapa de validação e comparação de resultados obtidos no reconhecimento do padrão maligno e benigno para as diferentes métricas de extração de características. O melhor resultado foi obtido pelo o índice de diversidade de Gleason utilizando a abordagem GLCM com acurácia e sensibilidade de 77%, e especificidade de 76%.

Palavras-chave: Mamografia . Textura . Imagens Medicas.

ABSTRACT

This work aims to investigate the application of texture analysis and pattern recognition techniques for the diagnosis of breast cancer, whose goal is to give greater support to the expert diagnosis of breast cancer, as we seek to use only to characterize the texture malignant and benign pattern of masses instead of contour features, since these features are not always clear in the images, since it may be overlapping as found, for example, to the masses and calcifications which have no lesions well defined boundary, from unauthorized viewing thereof and contributing to a larger number of negative biopsies.

This work proposes to study texture analysis techniques, such as: Local Binary Pattern and Gleason Diversity Indices and Menhinick, as it is believed that such techniques can produce good texture characteristics that discriminate against regions of masses in mammograms scanned between benign and malignant, since the success of the classification stage depends much on the characteristics generated. The produced features will be submitted as input for the classification process supervised using SVM and strategies combinations Ensemble models. In both cases, part of the samples will be used for classifier training step. This step creates a pattern on the extracted measures. The other part, totally unknown to the training stage, is used for testing and validation of the results. Finally follows the validation step and comparison of results in the recognition of malignant and benign pattern for the different feature extraction metrics. The best result was obtained by the Gleason diversity index using the GLCM approach with accuracy and sensitivity of 77% and specificity of 76%.

Keywords: Mammography. Texture. Medical Images.

SUMÁRIO

1	INTRODUÇÃO.....	14
1.2	Objetivos.....	15
1.2.1	Objetivos Específicos.....	15
1.3	Trabalhos Relacionados.....	15
1.4	Organização do trabalho.....	16
2	FUNDAMENTAÇÃO TEÓRICA.....	17
2.1	Processamento de Imagens Digitais.....	17
2.2	Quantização.....	18
2.3	Realce de Imagens.....	19
2.4	Matriz de Co-Ocorrência de Níveis de Cinza– GLCM.....	19
2.5	Matriz de Comprimento de Corrida de Cinza– GLRLM.....	21
2.6	Matriz de Comprimento de Lacuna de Cinza– GLGLM.....	23
2.7	Análise de Textura.....	24
2.8	Local Binary Pattern (LBP).....	25
2.9	Índice de Diversidade.....	25
2.9.1	Índice de Diversidade de Gleason.....	26
2.9.2	Índice de Diversidade de Menhinick.....	26
2.10	Reconhecimento de Padrões.....	27
2.10.1	SVM (Support Vector Machine).....	27
2.10.2	Adaboost.....	30
2.11	Validação de Resultados.....	31
3	METODOLOGIA.....	33
3.1	Aquisição de Imagens.....	33
3.2	Pré-Processamento.....	34
3.3	Extração de Características.....	35
3.3.1	Índice de Diversidade Ecológica.....	35
3.4	Reconhecimento de Padrões.....	38
3.4.1	SVM (Support Vector Machine).....	38
3.4.2	Adaboost.....	39
3.5	Validação dos Resultados.....	39
4	RESULTADOS E DISCUSSÃO.....	40
4.1	Índices de Diversidade.....	40

5	CONCLUSÃO.....	44
	REFERÊNCIAS BIBLIOGRÁFICAS.....	45

1 INTRODUÇÃO

A incidência do câncer de mama aumenta 1% ao ano. Segundo o INCA (Instituto Nacional de Câncer) (INCA, 2014), a previsão para 2014 era de 57.120 novos casos desse tipo de câncer e com um risco estimado de 56,09 casos a cada 100 mil mulheres.

O câncer de mama de acordo com (INCA, 2014) é o tipo de câncer que mais acomete as mulheres em todo mundo tanto em países em desenvolvimento quanto em países desenvolvidos. A idade é o principal fator de risco para o câncer de mama. As taxas de incidência aumentam rapidamente até os 50 anos e, posteriormente, esse aumento ocorre de forma mais lenta. Contudo, outros fatores de riscos já estão bem estabelecidos como: aqueles que são relacionados à vida produtiva da mulher (idade da primeira gestação acima do 30 anos, anticoncepcionais orais, menopausa tardia e terapia de reposição hormonal), histórico familiar de câncer de mama e alta densidade do tecido mamário (razão entre o tecido glandular e o tecido adiposo da mama). Além desses fatores de risco, a exposição à radiação ionizante, mesmo em baixas doses, também é considerada um fator de risco, particularmente durante a puberdade, segundo mostram alguns estudos.

A prevenção primária dessa doença ainda não é totalmente possível em decorrência da variação de fatores e das características genéticas que estão envolvidas na sua etiologia. Nova característica de rastreamento factível para países com dificuldades orçamentárias tem sido estudada, e, até o momento, a mamografia, para mulheres com idade entre 50 e 69 anos e o exame clínico das mamas anualmente a partir dos 40 anos, é recomendada como método efetivo para a detecção precoce. A amamentação, a prática de atividade física e a alimentação saudável com a manutenção do peso corporal estão associadas a um menor risco de desenvolver esse câncer (INCA, 2014).

Apesar de ser considerado um câncer de relativamente bom prognóstico se diagnosticado e tratado oportunamente, as taxas de mortalidade por essa doença continuam elevadas, provavelmente porque ela é diagnosticada em estágios avançados. A sobrevivência média após cinco anos na população de países desenvolvidos tem apresentado um pequeno aumento, cerca de 85%. Entretanto, em países em desenvolvimento, a sobrevivência é próxima a 60%. No Maranhão, para 2014, esse câncer tem uma taxa de 13,97 para cada 100 mil mulheres.

O diagnóstico do câncer de mama é feito principalmente através da mamografia, que é um exame de raios-x da mama e que permite a identificação de lesões da ordem de milímetros. A mamografia é atualmente uma das melhores técnicas de detecção precoce de

lesões não palpáveis na mama (ACS, 2014). Consiste em um exame de raios-X da mama, no qual o resultado é a produção, em uma folha de filme, de uma imagem em tons de cinza.

O diagnóstico baseado em mamografia é uma etapa feita por um radiologista, que ler e interpreta a imagem, assim diferentes especialistas podem fornecer interpretações diferentes para um mesmo exame. Além disso, a interpretação é uma tarefa repetitiva, requerendo um grande nível de atenção sobre os detalhes presentes na imagem.

Com as crescentes taxas de incidência e o risco de mortalidade se o mal não for detectado com antecedência, vários métodos foram e estão sendo pesquisados para a prevenção e detecção do câncer de mama. A computação gráfica e o processamento de imagens, assunto tratado aqui, são bastante utilizados na aplicação de imagens médicas.

1.2 Objetivos

Investigar a aplicação de técnicas de extração de características de textura e reconhecimento de padrões para caracterizar o padrão maligno e benigno de imagens de massas em mamografias digitalizadas cujo objetivo é dar ao especialista um maior suporte ao diagnóstico do câncer de mama, uma vez que busca-se utilizar somente a textura para caracterizar o padrão maligno e benigno ao invés das características do contorno das massas, já que tais características nem sempre são nítidas nas imagens.

1.2.1 Objetivos Específicos

- Estudar a viabilidade da utilização da técnica *Local Binary Pattern* (LBP) para extrair características de textura das massas;
- Estudar a viabilidade da utilização de técnicas de índice de diversidade, em especial, os índices de Gleason e Menhinick, para caracterizar a textura das massas;
- Estudar e aplicar a técnicas de reconhecimento de padrões: máquinas de vetor de suporte e estratégias de modelo Ensemble para testar as características produzidas no tocante ao poder de discriminação das classes malignas e benignas das massas.

1.3 Trabalhos Relacionados

Existe uma série de linhas de pesquisa relacionadas às diferentes técnicas aplicadas na metodologia computacional desse trabalho.

Em (ROCHA, 2014) é desenvolvida uma metodologia para discriminar padrões de malignidade e benignidade de massas em imagens de mamografias, utilizando abordagens

estruturais e estatísticas para análise de textura. As técnicas utilizadas foram: Índices de Diversidade Shannon, Mcintosh, Simpson, Gleason e Menhinick, *Local Binary Pattern*, Função K de Ripley e Matrizes de Co-Ocorrência. A textura extraída é classificada usando a Máquina de Vetores de Suporte. O melhor resultado foi obtido usando a função K de Ripley com 92,20% de acurácia, 92,96% de sensibilidade e 91,26% de especificidade.

Em (CARVALHO, 2012) é utilizado o Índice de Diversidade de Mcintosh e a Máquina de Vetores de Suporte para discriminação e classificação de regiões extraídas de mamografias em massas e não-massas. Para calcular os índices são utilizadas quatro abordagens: o Histograma, a Matriz de Co-Ocorrência de Níveis de Cinza, a Matriz de Comprimento de Lacuna de Cinza e a Matriz de Comprimento de Corrida de Cinza. O melhor resultado apresentou uma acurácia de 93,68%.

Em (LIU, 2011) e (NANNI, 2012) utiliza-se o *Local Binary Pattern* para extrair características de textura e assim discriminar massas quanto ao seu caráter maligno e benigno. Já em (MASCARO, 2009), o LBP é utilizado para segmentar imagens de mamografias para detecção de massas.

A partir da análise dos trabalhos relacionados verifica-se que a análise da textura como parâmetro para extração de características e reconhecimento de padrões apresentam resultados promissores na detecção de câncer de mama, baseado em mamografias.

1.4 Organização do trabalho

Este trabalho apresenta a seguinte organização:

No Capítulo 2, A Fundamentação Teórica traz informações importantes para o contexto e entendimento do trabalho, tais como explicações sobre conceitos de Processamento de Imagens, *Local Binary Pattern* e Índices de Diversidade e Reconhecimento de Padrões.

No Capítulo 3, Metodologia, será explicada a metodologia utilizada como ponto de partida para o desenvolvimento desse trabalho.

No Capítulo 4, Resultados, serão apresentados os resultados alcançados na aplicação da metodologia, junto com uma análise de seu desempenho.

No Capítulo 5, Conclusão, apresenta-se a conclusão do trabalho. Nela está contida uma retrospectiva do que foi falado na monografia como um todo e uma avaliação dos resultados obtidos.

2 FUNDAMENTAÇÃO TEÓRICA

Nessa seção, abordaremos os conceitos necessários para entender o funcionamento da metodologia computacional proposta.

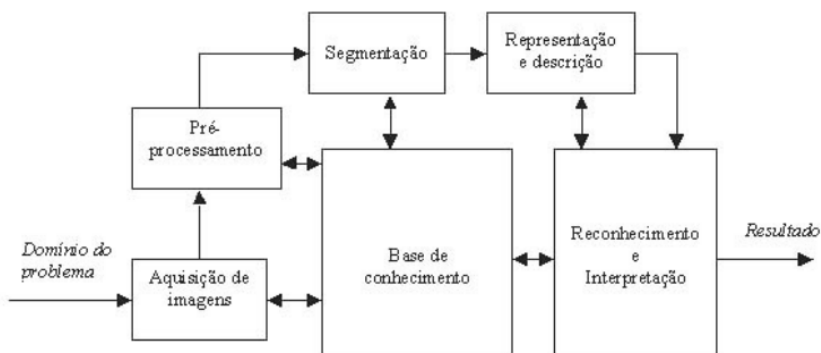
2.1 Processamento de Imagens Digitais

Uma imagem pode ser definida como uma função bidimensional, $f(x, y)$, onde x e y são as coordenadas espaciais, e a amplitude de f em qualquer par de coordenadas (x, y) é a intensidade ou o chamado nível de cinza da imagem em um ponto (GONZALEZ, 2002).

Assim, o processamento de imagens digitais é composto por um conjunto de técnicas que englobam desde a aquisição até seu reconhecimento e interpretação.

A área de processamento de imagens vem sendo objeto de crescente interesse por permitir viabilizar grande número de aplicações em duas categorias bem distintas: (1) o aprimoramento de informações pictóricas para interpretação humana; e (2) a análise automática por computador de informações extraídas de uma cena (MARQUES, 1999).

Figura 1. Etapas do processamento de imagens.



Fonte: (GONZALEZ, 1992).

O primeiro passo no processamento de imagens é adquirir uma imagem digital. Esta etapa tem como função converter uma imagem em uma representação numérica adequada para o processamento digital subsequente. Este bloco compreende dois elementos principais. O primeiro é um dispositivo físico sensível a uma faixa de energia no espectro eletromagnético (como raio-X, ultravioleta, espectro visível ou raios infravermelhos), que produz na saída um sinal elétrico proporcional ao nível de energia detectado. O segundo - o

digitalizador propriamente dito - converte o sinal elétrico analógico em informação digital (MARQUES, 1999).

O pré-processamento tem como função chave melhorar a imagem de forma a aumentar as chances para o sucesso dos processos seguintes. Tipicamente envolve técnicas para o realce de contrastes, remoção de ruído e isolamento de regiões cuja textura indique a probabilidade de informação alfanumérica.

O próximo estágio trata da segmentação que divide uma imagem de entrada em partes ou objetos constituintes. Em geral, a segmentação automática é uma das tarefas mais difíceis no processamento de imagens digitais.

Por um lado, um procedimento de segmentação robusto favorece substancialmente a solução bem sucedida de um problema de imageamento. Por outro lado, algoritmos de segmentação fracos ou erráticos quase sempre asseveram falha no processamento. No caso de reconhecimento de caracteres, o papel básico da segmentação é extrair caracteres individuais e palavras do fundo da imagem.

A saída do estágio de segmentação é constituída tipicamente por dados em forma de pixels, correspondendo tanto à fronteira de uma região como a todos os pontos dentro da mesma.

O processo de descrição, também chamado de seleção de características, procura extrair características que resultem em alguma informação quantitativa de interesse ou que sejam básicas para discriminação entre classes de objetos. Em se tratando de reconhecimento de caracteres, descritores tais como buracos e concavidades são características poderosas que auxiliam na diferenciação entre uma parte do alfabeto e outra.

O último estágio envolve reconhecimento e interpretação. Reconhecimento é o processo que atribui um rótulo a um objeto, baseado na informação fornecida pelo seu descritor. A interpretação envolve a atribuição de significado a um conjunto de objetos reconhecidos.

Nas próximas seções serão apresentados os conceitos e os métodos de Processamento de Imagens utilizadas no desenvolvimento deste trabalho.

2.2 Quantização

A quantização é um processo que procura obter a representação de uma imagem com L níveis de cinza para cada pixel, com $L = 2^b$, sendo b o número de *bits* usados para armazenar o valor do *pixel*. Assim, dada uma imagem com L níveis de cinza, se houver

necessidade de quantizá-la para L' níveis de cinza, onde, $L' < L$ podemos usar a quantização uniforme, que consiste em dividir a escala de cinza da imagem em intervalos iguais, nos quais cada intervalo é mapeado para um valor cinza na imagem quantizada, de modo que a escala de cinza da imagem quantizada é dada por $[0, L' - 1]$ (GONZALEZ, 1992).

$$q(i, j) = (2^b - 1) \frac{p(i, j) - I_{min}}{I_{max} - I_{min}} \quad (1)$$

onde $q(i, j)$ é o nível de cinza do *pixel* (i, j) da nova imagem (quantizada), $p(i, j)$ é o nível de cinza do *pixel* (i, j) da imagem original, $[I_{max} - I_{min}]$ são os limites inferior e superior da escala de cinza da imagem original e b é o número de *bits* necessários para armazenar cada *pixel* da imagem quantizada.

2.3 Realce de Imagens

O realce de imagens é o processo de manipular uma imagem de forma que o resultado seja mais adequado do que o original para uma aplicação específica. Assim, não existe uma técnica de realce que possa ser aplicada a qualquer categoria de problema, uma vez que as técnicas de realce são orientadas ao problema (GONZALEZ, 2010).

No caso das aplicações médicas, em especial as imagens de mamografia, algumas técnicas de realce merecem destaque, entre elas está o realce logarítmico.

Esta técnica aumenta o contraste em valores de cinza baixos. Equivale a uma curva logarítmica. É definida como:

$$g_t(l, p) = G \log_{10}(g(l, p) + 1); G = \frac{max}{\log_{10} max} \quad (2)$$

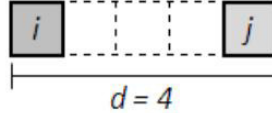
onde $g_t(l, p)$ representa o valor do nível de cinza do ponto (l, p) , $g(l, p)$ é o valor original do *pixel*, G é gerado a partir do valor de max , que é o limite máximo de nível de cinza na imagem.

2.4 Matriz de Co-Ocorrência de Níveis de Cinza– GLCM

De acordo com (PEDRINI, 2008), dado um relacionamento espacial entre os *pixels* componentes de uma textura, os elementos da matriz de co-ocorrência de níveis de cinza (GLCM, do inglês *Gray Level Co-Occurrence Matrix*) descrevem a frequência com que ocorrem as transições de nível de cinza entre pares de *pixels*. Efetuando-se variações na

relação espacial, por meio de alterações na orientação e distância entre as coordenadas dos *pixels*, podem ser obtidas diversas matrizes de co-ocorrência, a partir das quais são extraídas medidas utilizadas para análise de texturas (HARALICK, 1973).

Figura 2. Exemplo de uma co-ocorrência dos níveis de cinza i e j , com vizinhança $d = 4$, alinhados na horizontal ($\theta = 0$).



Fonte: (CARVALHO, 2012)

Dada uma imagem S , com níveis de cinza no intervalo $[0, L - 1]$, cada célula (i, j) da matriz de co-ocorrência, com $0 \leq i \leq L - 1$ e $0 \leq j \leq L - 1$ funciona como um contador e armazena a frequência, denotada por $P(i, j, d, \theta)$, com que dois *pixels* ocorrem na imagem, separados por uma distância d , sob uma direção θ , um com a cor i e outro com a cor j . O cálculo do elemento da matriz de co-ocorrência, para as direções $0^\circ, 45^\circ, 90^\circ$ e 135° , é descrito através de 4 equações (HARALICK, 1973):

$$P(i, j, d, 0^\circ) = \#\{(k, l), (m, n) \mid k - m = 0, |l - n| = d, f(k, l) = i, f(m, n) = j\} \quad (3)$$

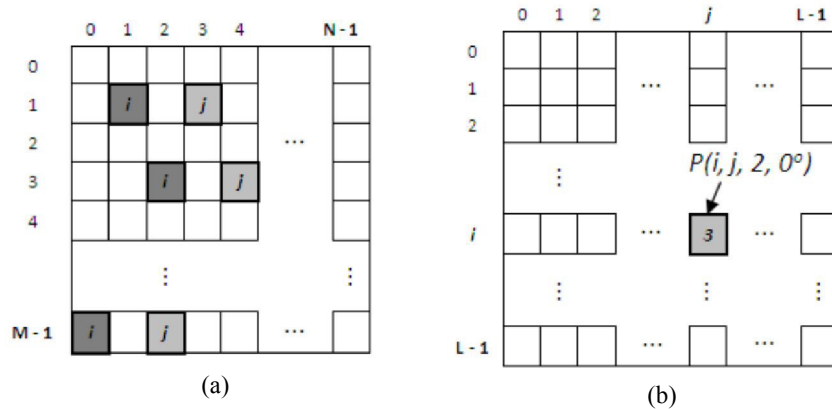
$$P(i, j, d, 45^\circ) = \#\{(k, l), (m, n) \mid k - m = d, l - n = -d, f(k, l) = i, f(m, n) = j\} \quad (4)$$

$$P(i, j, d, 90^\circ) = \#\{(k, l), (m, n) \mid |k - m| = d, l - n = 0, f(k, l) = i, f(m, n) = j\} \quad (5)$$

$$P(i, j, d, 135^\circ) = \#\{(k, l), (m, n) \mid k - m = d, l - n = d, f(k, l) = i, f(m, n) = j\} \quad (6)$$

onde “#” denota o número de pares $((k, l), (m, n))$ do conjunto, e $f(x, y)$ denota a função de cinza no *pixel* (x, y) .

Figura 3. (a) Imagem de $M \times N$ pixels. (b) Matriz de Co-Ocorrência de Níveis de Cinza da imagem ($d = 2, \theta = 0$).



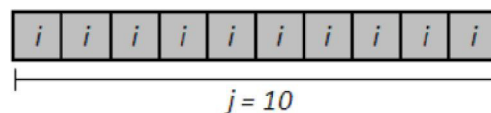
Fonte: (CARVALHO, 2012)

A Figura 3(b) ilustra a estrutura da GLCM, construída a partir da imagem da Figura 3(a). O tamanho da matriz é $L \times L$, sendo L a quantidade máxima de níveis de cinza que a imagem pode apresentar. Na Figura 3(a), por exemplo, há 3 pares de *pixels*, com vizinhança 2 e alinhamento na horizontal, onde o primeiro *pixel* tem intensidade i e o segundo tem intensidade j . Assim, a célula (i, j) da GLCM registra a frequência $P(i, j, 2, 0^\circ) = 3$.

2.5 Matriz de Comprimento de Corrida de Cinza– GLRLM

Dada uma imagem, define-se que um conjunto composto de *pixels* consecutivos, apresentando o mesmo nível de cinza e sendo colineares em uma dada direção, representa uma corrida de cinza. O número de *pixels* contidos nesse conjunto denota o comprimento da corrida.

Figura 4. Exemplo de uma corrida de nível de cinza i , de comprimento 10 e direção horizontal.



Fonte: (CARVALHO, 2012)

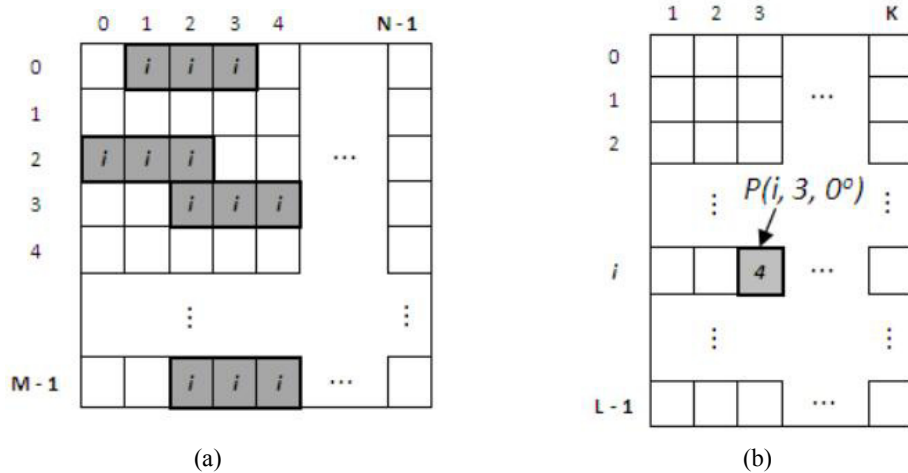
Com o objetivo de sintetizar as informações obtidas a partir dessas corridas, são criadas matrizes de comprimentos de corrida de cinza (GLRLM, do inglês *Gray Level Run*

Length Matrix), nas quais cada elemento, representado por $P(i, j, \theta)$, contém o número de corridas com tamanho j (comprimento), tendo i como o nível de cinza de seus pixels, e o parâmetro θ como a orientação do segmento de reta formado pelos *pixels*. A partir da GLRLM podem ser extraídas medidas usadas para análise de textura (GALLOWAY, 1975). O cálculo do elemento da GLRLM (BEBIS, 2006) é definido como a seguir:

$$P(i, j, \theta) = \text{CARD}[\{(m, n) | f(m, n) = i, \tau(i, \theta) = j\}] \quad (7)$$

onde $f(m, n)$ denota a função nível de cinza no *pixel* (m, n) , e $\tau(i, \theta)$ é o comprimento da corrida de nível de cinza i e direção θ , e CARD significa a cardinalidade (número de elementos) do conjunto. Os valores de θ adotados são $0^\circ, 45^\circ, 90^\circ$ e 135° .

Figura 5. (a) Imagem de $M \times N$ pixels. (b) Matriz de Comprimento de Corrida de Níveis de Cinza da imagem ($\theta = 0^\circ$).



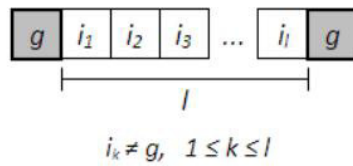
Fonte: (CARVALHO, 2012)

A Figura 5(b) ilustra a estrutura da GLRLM, construída a partir da imagem da Figura 5(a). O tamanho da matriz é $L \times K$, sendo L a quantidade máxima de níveis de cinza que a imagem pode apresentar e K o maior comprimento de corrida presente na imagem na direção θ . Na Figura 5(a), por exemplo, há 4 corridas de *pixels* com nível de cinza i , comprimento 3 e direção horizontal. Assim, a célula $(i, 3)$ da GLRLM registra a frequência $P(i, 3, 0^\circ) = 4$.

2.6 Matriz de Comprimento de Lacuna de Cinza– GLGLM

Dada uma imagem, define-se que uma lacuna para o nível de cinza g ocorre quando g é encontrado apenas no início e no fim de um conjunto de pixels consecutivos e colineares, enquanto todos os valores de *pixels* entre g , estão acima ou abaixo de g (Figura 6). O comprimento da lacuna é a distância entre estes dois *pixels* menos um, de modo que, dois *pixels* vizinhos adjacentes com nível de cinza idêntico têm comprimento de lacuna zero. No caso em que nenhum *pixel* com nível de cinza g é encontrado ao longo da direção de busca, o comprimento da lacuna é considerado como infinito, sendo omitido (XINLI, 1994).

Figura 6. Lacuna de nível de cinza g , de comprimento l e direção horizontal.



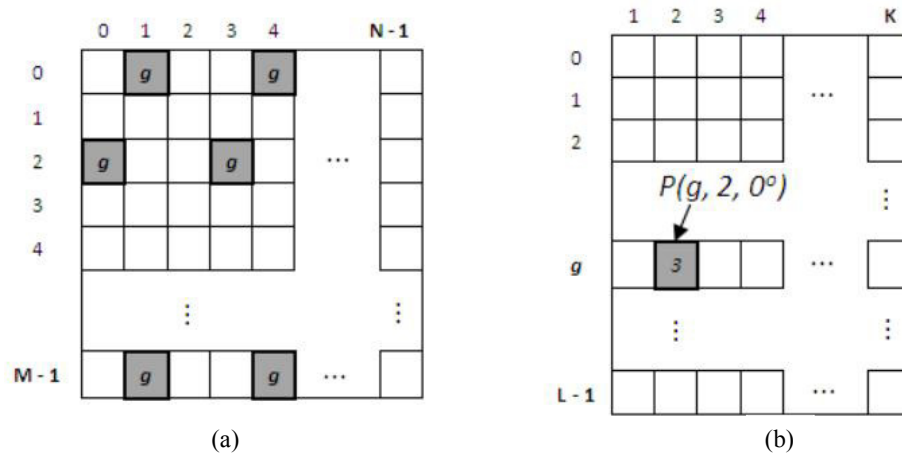
Fonte: (CARVALHO, 2012)

A matriz de comprimento de lacuna de nível de cinza (GLGLM, do inglês *Gray Level Gap Length Matrix*) é uma matriz estatística de ordem superior, na qual cada elemento (g, l) armazena a frequência denotada por $P(g, l, \theta)$, com que lacunas de nível de cinza g , tamanho l , e inclinação θ ocorrem na imagem (Figura 7). O elemento da GLGLM (XINLI, 1994), na direção θ , é definido como:

$$\begin{aligned}
 P(g, l, \theta) &= \#\{(i, j) | f(i, j) = g, \\
 &f(i + x, j + y) = g, \\
 &f(i + u, j + v) \neq g, \\
 &x = (l + 1) \cdot \cos \theta, \\
 &y = (l + 1) \cdot \sin \theta, \\
 &u < x, v < y\}
 \end{aligned} \tag{8}$$

onde “#” denota o número de elementos do conjunto, e $f(i, j)$ denota a função nível de cinza no *pixel* (i, j) .

Figura 7. (a) Imagem de $M \times N$ pixels. (b) Matriz de Comprimento de Lacuna de Níveis de Cinza da imagem ($\theta = 0$).



Fonte: (CARVALHO, 2012)

A Figura 7(b) ilustra a estrutura da GLGLM, construída a partir da imagem da Figura 7(a). O tamanho da matriz é $L \times K$, sendo L a quantidade máxima de níveis de cinza que a imagem pode apresentar e K o maior comprimento de corrida presente na imagem na direção θ . Na Figura 7(a), por exemplo, há 3 lacunas de *pixels* com nível de cinza g , comprimento 2 e inclinação horizontal. Assim, a célula $(g, 2)$ da GLGLM registra a frequência $P(g, 2, 0^\circ) = 3$.

2.7 Análise de Textura

A textura pode ser compreendida como um descritor que fornece medidas de uniformidade, densidade, aspereza, regularidade, intensidade, entre outras características da imagem. Em (BRAZ, 2008), a textura é definida como uma propriedade importante de percepção de região e superfícies, contendo informações sobre a distribuição espacial das variações de tonalidade locais em valores de *pixels* que se repetem de maneira regular ou aleatória ao longo do objeto ou imagem.

Em processamento de imagens, utilizam-se três métodos principais para descrever a textura de uma região: a estatística, a estrutural e a espectral.

Neste trabalho utilizamos a abordagem estatística, que caracteriza a textura como suave, rugosa, granulada, etc, e a abordagem estrutural, que trata a textura como subpadrões espaciais na imagem. Para extrair as medidas que serviram para caracterizar a textura usamos o *Local Binary Pattern* (LBP) e os índices de diversidade de Gleason e Menhinick.

2.8 Local Binary Pattern (LBP)

O *Local Binary Pattern* é um operador não paramétrico para descrever a estrutura espacial local da imagem, mostrando alta capacidade de distinguir características de textura (OJALA, 1996).

A função abaixo descreve o cálculo do LBP,

$$LBP(x_c, y_c) = \sum_{n=0}^{n-1} S(i_n - i_c) 2^n \quad (9)$$

onde n é o número de vizinhos do *pixel* central x_c, y_c considerados no cálculo, i_c é o valor de nível de cinza do *pixel* central x_c, y_c , i_n é o valor de nível de cinza de cada *pixel* vizinho e $S(x)$ uma função que devolve 1 se $x \geq 0$ e 0, caso contrário.

Figura 8. Exemplo do cálculo do LBP.

6	5	2	1	0	0	1	2	4	1	0	0
7	6	1	1		0	8		16	8		0
9	3	7	1	0	1	32	64	128	32	0	128
a)			b)			c)			d)		

Fonte: (ROCHA, 2014)

Um exemplo do cálculo do LBP é apresentado na Figura 8. Dada uma janela de tamanho 3 x 3 (Figura 8(a)) centrada em um *pixel*, é feita a subtração dos valores dos níveis de cinza dos *pixels* vizinhos (um por vez) com o valor do nível de cinza do *pixel* central, formando uma matriz binária composta pelo correspondente valor 0 ou 1, dependendo do resultado da diferença dos *pixels* analisados (Figura 8(b)). Estes valores da matriz binária são multiplicados pelo respectivo valor da matriz de pesos (Figura 8(c)). O LBP é o resultado da soma de todos os valores resultantes da multiplicação (Figura 8(d)). No exemplo em questão, o LBP é 169.

2.9 Índice de Diversidade

O termo diversidade é definido como a variedade de espécies presentes em uma comunidade, habitat ou região. Uma comunidade é definida como um conjunto de espécies que ocorrem em um determinado lugar e tempo. Os índices de diversidade são usados para

representar a composição de uma comunidade, possibilitando o dimensionamento de sua riqueza, igualdade e a diversidade das espécies nos diferentes ambientes estudados (MAGURRAN, 2004).

Ainda no âmbito ecológico, por definição, a diversidade envolve dois parâmetros: riqueza, que corresponde a quantidade de espécies; e abundância relativa, que é a quantidade de indivíduos de determinada espécie, que ocorre em um local ou amostra. Assim, segundo McIntosh, comunidades com a mesma riqueza podem diferir em diversidade dependendo da distribuição de indivíduos entre as espécies.

Uma medida de diversidade é um parâmetro extremamente reducionista que objetiva expressar toda a complexidade estrutural de uma comunidade ecológica através de um único número. Assim, é vantajoso o fato do índice de diversidade utilizar um único número para representar uma determinada situação, já que facilita a comparação em experimentação, e também possibilita a elucidação de mudanças que ocorrem nas comunidades relacionadas (SANTOS, 2009).

2.9.1 Índice de Diversidade de Gleason

O índice de diversidade de Gleason é um índice simples, pois considera somente o número de espécie (s) e o logaritmo (base 10 ou natural) do número total de indivíduos (BROWER, 1997).

$$D_g = \frac{s}{\log N} \quad (10)$$

Classificado como um índice de riqueza, pois expressa a riqueza (s) como o número simples de espécies, ocorrendo a alteração deste valor em função do tamanho da amostra.

2.9.2 Índice de Diversidade de Menhinick

O índice de diversidade de Menhinick (1964) também é considerado simples, pois considera somente o número de espécies (s) e a raiz quadrada do número total de indivíduos, sendo calculado pela equação:

$$D_b = \frac{s}{\sqrt{N}} \quad (11)$$

onde s é o número de espécies amostradas e N é o número total de indivíduos em todas as espécies. Assim como o índice de Gleason, o índice de Menhinick também é classificado como um índice de riqueza de espécies.

2.10 Reconhecimento de Padrões

Em (LOONEY, 1997), um padrão é definido como tudo aquilo para o qual existe uma entidade nomeável representante, geralmente criada através do conhecimento cultural humano. O reconhecimento de padrões visa determinar um mapeamento que relacione as propriedades extraídas de amostras com um conjunto de rótulos (entidade nomeável representante), apresentando a restrição de que amostras com características semelhantes devem ser mapeadas ao mesmo rótulo. Os algoritmos que estabelecem este mapeamento são denotados como algoritmos de classificação ou classificadores (PEDRINI, 2008).

O processo de classificação pode ser feito de duas formas, supervisionada e não-supervisionada. A classificação supervisionada ocorre quando o classificador considera classes pré-definidas e uma etapa de treinamento deve ser executada para que os parâmetros que caracterizam cada classe sejam obtidos. Na classificação não-supervisionada não se dispõe de parâmetros ou informações coletadas previamente à aplicação do algoritmo de classificação, e todas as informações devem ser obtidas a partir das próprias amostras a serem rotuladas (PEDRINI, 2008).

Neste trabalho utilizam-se dois processos de classificação supervisionada, o SVM (*Support Vector Machine*) e o *Adaboost* para realizarem o reconhecimento de padrões de tecidos da mama (massas), de modo a determinar sua natureza maligna ou benigna.

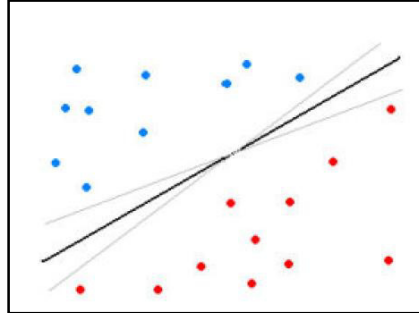
2.10.1 SVM (Support Vector Machine)

Support Vector Machine (SVM) é uma técnica de aprendizagem supervisionada usada para estimar uma função que classifique dados de entrada em duas classes. O princípio básico é a construção de um hiperplano como superfície de decisão, cuja margem de separação entre as classes seja máxima. Por hiperplano entende-se uma superfície de separação de duas regiões em um espaço multidimensional, no qual o número de dimensões pode ser até infinito (VAPNIK, 2008).

A Figura 9 mostra em duas dimensões, para melhor visualização, hiperplanos de separação entre duas classes linearmente separáveis. O hiperplano ótimo (linha mais escura),

não somente separa as duas classes, mas mantém a maior distância possível com relação aos pontos da amostra. Há casos em que podem existir vários possíveis hiperplanos de separação, mas o SVM busca apenas encontrar o que maximize a margem entre os exemplos de treinamento.

Figura 9. Separação entre duas classes através de hiperplanos.



Fonte: (NASCIMENTO, 2012).

Seja o conjunto de amostras de treinamento (x_i, y_i) , sendo $x_i \in \mathbb{R}^n$ o vetor de entrada, y_i a classificação correta das amostras e $i=1,2,\dots,n$ o índice de cada ponto amostral. O objetivo da classificação é estimar a função $f: \mathbb{R}^n \rightarrow \{\pm 1\}$ que separe corretamente os exemplos de teste em classes distintas.

A etapa de treinamento estima a função $f(x) = (w \cdot x) + b$, procurando valores tais que a seguinte relação seja satisfeita:

$$y_i((w \cdot x_i) + b) \geq 1 \quad (12)$$

Sendo w o vetor normal ao hiperplano de decisão e b o corte ou distância da função f em relação à origem. Os valores ótimos de w e b serão encontrados de acordo com a restrição dada pela Equação 12 ao minimizar a seguinte equação:

$$\phi(w) = \frac{w^2}{2} \quad (13)$$

O SVM ainda possibilita encontrar um hiperplano que minimize a ocorrência de erros de classificação nos casos em que uma perfeita separação entre as duas classes não seja possível. Isso graças à inclusão de variáveis de folga, que permitem que as restrições presentes na Equação 12 sejam quebradas.

O problema de otimização passa a ser então a minimização da Equação 14, de acordo com a restrição imposta na Equação 12. C é um parâmetro de treinamento que estabelece um

equilíbrio entre a complexidade do modelo e o erro de treinamento e deve ser selecionado pelo usuário.

$$\phi(w, \xi) = \frac{w^2}{2} + C \sum_{i=1}^N \xi_i \quad (14)$$

para

$$y_i((w \cdot x_i) + b) + \xi_i \geq 1 \quad (15)$$

Através da teoria dos multiplicadores de Lagrange, chega-se à Equação 16. O objetivo então passa a ser encontrar os multiplicadores de Lagrange a_i ótimos que satisfaçam a Equação 17 (CHAVES, 2006).

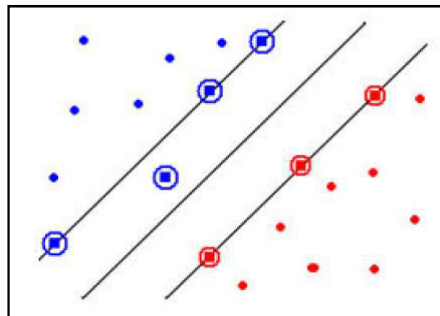
$$L(a) = \sum_{i=1}^N a_i - \frac{1}{2} \sum_{i,j=1}^N a_i a_j y_i y_j (x_i, x_j) \quad (16)$$

$$\sum_{i=1}^N a_i y_i = 0, \quad 0 \leq a_i \leq C \quad (17)$$

Apenas os pontos onde a restrição dada pela Equação 12 é exatamente igual à unidade têm correspondentes $a_i \neq 0$. Esses pontos são chamados de vetores de suporte, pois se localizam geometricamente sobre as margens. Tais pontos têm fundamental importância na definição do hiperplano ótimo, pois os mesmos delimitam a margem do conjunto de treinamento. A Figura 10 destaca os pontos que representam os vetores de suporte.

Os pontos além da margem não influenciam decisivamente na determinação do hiperplano, enquanto que os vetores de suporte, por terem pesos não nulos, são decisivos.

Figura 10. Vetores de Suporte (destacado por círculos).



Fonte: (NASCIMENTO, 2012).

Para que o SVM possa classificar amostras que não são linearmente separáveis, é necessária uma transformação não-linear que transforme o espaço de entrada (dados) para um novo espaço (espaço de características).

Esse espaço deve apresentar dimensão suficientemente grande, e através dele, a amostra pode ser linearmente separável. Dessa maneira, o hiperplano de separação é definido como uma função linear de vetores retirados do espaço de características ao invés do espaço de entrada original. Essa construção depende do cálculo de uma função K de núcleo de um produto interno (HAYKIN, 2001). A função k pode realizar o mapeamento das amostras para um espaço de dimensão muito elevada sem aumentar a complexidade dos cálculos.

A Equação 18 mostra o resultado da Equação 16 com a utilização de um núcleo k .

$$L(a) = \sum_{i=1}^N a_i - \frac{1}{2} \sum_{i,j=1}^N a_i a_j y_i y_j K(x_i, x_j) \quad (18)$$

Uma importante família de funções de núcleo é a função de base radial, muito utilizada em problemas de reconhecimento de padrões e também empregada neste projeto. A função de base radial é definida por:

$$K(x_i, y_j) = \exp(-\gamma \|x_i - y_j\|^2) \quad (19)$$

onde $\gamma = 1/\sigma^2$, sendo σ a variância.

2.10.2 Adaboost

O *Adaboost* é um algoritmo de aprendizado de máquina do tipo *boosting*. O *boosting* funciona em iterações. A cada iteração, um algoritmo-base é chamado para gerar um classificador simples, utilizando uma diferente versão do conjunto de dados de treinamento. As diferentes versões do conjunto de treinamento são obtidas através da variação do peso associado a cada um dos exemplos. Assim, temos diferentes versões ponderadas do conjunto de dados. Após um número determinado de iterações, o *boosting* combina os diversos classificadores parciais, gerando um classificador único, que, com o intuito de obter um melhor desempenho do que o do melhor classificador parcial (DUARTE, 2009).

O *Adaboost* chama um algoritmo-base em várias iterações t , onde $t \in [1 \dots T]$. Em cada iteração t , a distribuição de pesos do conjunto de treinamento é atualizada para a utilização pelo algoritmo-base. A atualização é realizada de forma a, relativamente, aumentar os pesos dos exemplos incorretamente classificados em confronto com os pesos dos exemplos corretamente classificados (DUARTE, 2009).

Algoritmo 1: Pseudocódigo do *Adaboost*.

-
- 1 Obtém N exemplos de imagens $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$, onde x corresponde a matriz de pixels de uma imagem, e $y = 0$, ou 1 para exemplos negativos e positivos, respectivamente.
 - 2 Inicializa os pesos $\omega_{1,i} = \frac{1}{2m}, \frac{1}{2l}$ para $y_i = 0$ e 1 respectivamente, onde m é o total de exemplos negativos e l o total de positivos.
 - 3 **para** $t=1, \dots, T$ **faça**
 - 4 Normalize os pesos: $\omega_{t,i} \leftarrow \frac{\omega_{t,i}}{\sum_{j=1}^N \omega_{t,j}}$, onde ω_t é uma distribuição de probabilidade.
 - 5 **para cada característica** j **faça**
 - 6 Treine um classificador h_j
 - 7 Avalie o erro de acordo com: $\omega_j, \epsilon_j = \sum_i \omega_i |h_j(x_i) - y_i|$.
 - 8 Escolha o classificador h_t com o menor erro ϵ_t .
 - 9 Atualize os pesos: $\omega_{t+1,i} = \omega_{t,i} \beta_t^{(1-e_i)}$, onde $e_i = 0$, se o exemplo x_i for classificado corretamente, $e_i = 1$ caso contrário, e $\beta_t = \frac{\epsilon_t}{1-\epsilon_t}$.
 - 10 O classificador forte é definido por: $H(x) = \begin{cases} 1, & \text{se } \sum_{t=1}^T \alpha_t h_t(x) \geq \frac{1}{2} \sum_{t=1}^T \alpha_t \\ 0, & \text{caso contrário.} \end{cases}$,
onde $\alpha_t = \log \frac{1}{\beta_t}$.
-

O *adaboost.M1* é uma adaptação do *boosting* de (FREUND, 1995), para problemas de classificação multi-classe. A principal diferença é a hipótese final H que, dada uma instancia x , tem como saída o valor y que maximiza a soma dos valores α_t que predizem uma classe (DUARTE, 2009).

$$H(x) = \underset{y \in Y}{\operatorname{argmax}} \left(\sum_{t=1}^T \alpha_t [h_t(x) = y] \right) \quad (20)$$

2.11 Validação de Resultados

A validação dos resultados produzidos é uma etapa que busca medir o desempenho da metodologia, calculando-se algumas estatísticas sobre os resultados dos testes e também porque a etapa de reconhecimento de padrões é um processo que resulta mais em probabilidade de se estar certo do que na certeza propriamente dita.

Na análise de imagens médicas, geralmente utiliza-se algumas estatísticas descritivas sobre os resultados dos testes para avaliar o desempenho do classificador, como sensibilidade (S), especificidade (E) e acurácia (A) (BLAND, 2000). Estas métricas são calculadas a partir de quatro situações possíveis em relação ao diagnóstico:

- VP – Verdadeiro Positivo: o teste é positivo e o paciente tem a doença;
- FP – Falso Positivo: o teste é positivo, mas o paciente não tem a doença;
- VN – Verdadeiro Negativo: o teste é negativo e o paciente não tem a doença;
- FN – Falso Negativo: o teste é negativo, mas o paciente tem a doença;

A acurácia corresponde a taxa de casos classificados corretamente sobre o numero total de casos:

$$A = \frac{VP + VN}{VP + FP + VN + FN} \quad (21)$$

A sensibilidade define a proporção de pessoas com a doença de interesse que têm o resultado do teste positivo. Indica quão bom é o teste para identificar indivíduos doentes:

$$S = \frac{VP}{VP + FN} \quad (22)$$

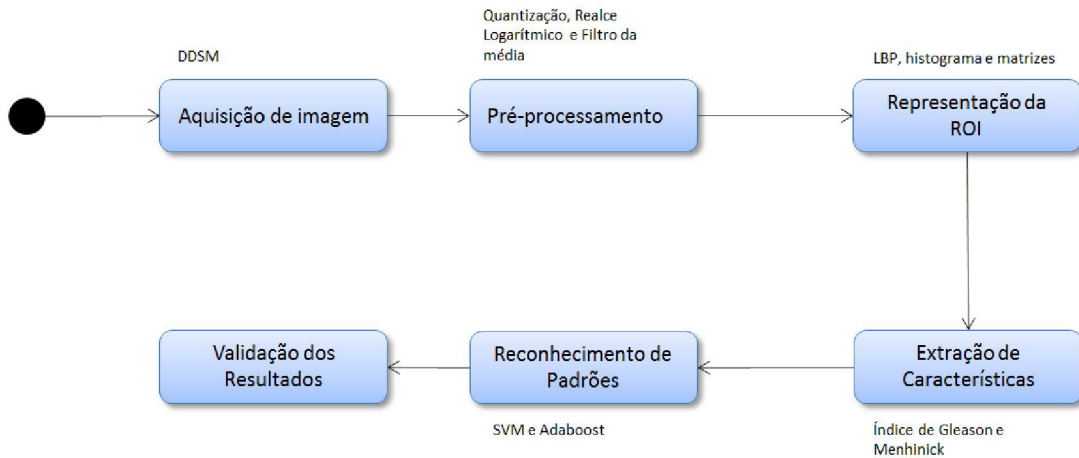
A especificidade define a proporção de pessoas sem a doença de interesse que tem o resultado do teste negativo. Indica quão bom é o teste para identificar indivíduos não doentes:

$$E = \frac{VN}{VN + FP} \quad (23)$$

3 METODOLOGIA

Este capítulo apresenta a metodologia proposta para a diferenciação dos padrões malignos e benignos de massas, a partir de imagens de mamografia. A Figura 11 apresenta as etapas da metodologia.

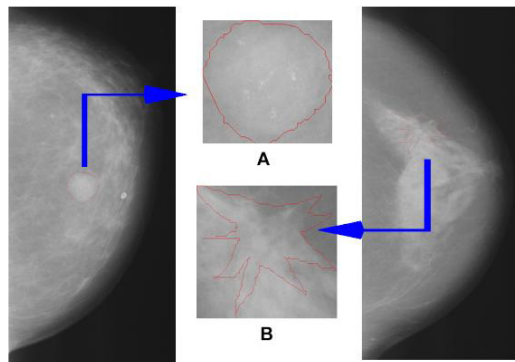
Figura 11. Etapas da metodologia proposta.



3.1 Aquisição de Imagens

Nesta etapa da metodologia são obtidas as amostras de mamografia que foram empregadas nos testes. Assim utilizou-se a base pública, de mamografias digitalizadas, DDSM (*Digital Database for Screening Mammography*), disponível na internet. A base é formada por 2620 exames de pacientes de diferentes origens étnicas e raciais. Cada exame contém duas imagens de cada mama, nas projeções médio-lateral oblíqua e crânio-caudal. Além disso, são disponibilizadas informações sobre a paciente, tal como a idade e a densidade da mama. Junto com as imagens que apresentam áreas suspeitas (massas) é fornecido um arquivo de descrição de lesão (overlay), contendo a quantidade de lesões presentes na mamografia, a localização da lesão, o tipo de lesão, o contorno da lesão e seu diagnóstico. O contorno da lesão está codificado em *chain code* (MORSE, 2000).

Figura 12. (a) Mamografia com nódulo benigno selecionado. (b) Mamografia com nódulo maligno selecionado.



Fonte: (BRAZ, 2008).

Como o objetivo desta monografia foi a caracterização da textura das massas usando LBP, índices de diversidade e posteriormente, sua classificação quanto à natureza maligna ou benigna, não foi utilizada a imagem completa da mamografia, partindo-se do pressuposto que as ROIs foram extraídas anteriormente. Assim como em (BRAZ, 2008) e (ROCHA, 2014), utilizaram-se as amostras selecionadas a partir das bounding boxes, obtendo-se somente as regiões que contem as massas, totalizando 3559 ROIs. Na etapa, de teste utilizou-se um subconjunto de 1155 ROIs, sendo 625 massas malignas e 530 benignas.

3.2 Pré-Processamento

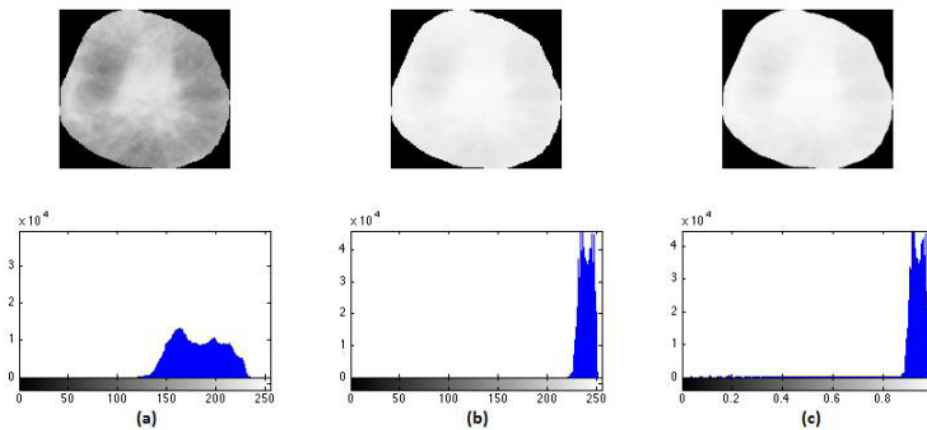
Após a aquisição das amostras, as mesmas foram submetidas a um pré-processamento que tem por objetivo dar um melhor contraste ao objeto de interesse em relação ao fundo da imagem e assim obter uma melhor descrição da textura das massas, para isso, utilizou-se o realce logarítmico seguido de uma suavização pelo filtro da média.

Como em imagens de mamografia os pixels com níveis de cinza mais escuros têm uma frequência menor, o realce logarítmico foi utilizado para dar maior relevância aos mesmos.

O realce logarítmico faz uso da definição de uma constante G , que é determinada a partir dos limites mínimo e máximo da imagem, e busca garantir que os novos valores estejam entre 0 e 255. Assim o valor utilizado neste trabalho para a constante G será igual a 105,98.

O filtro da média tem o objetivo de remover picos gerados após a aplicação do realce logarítmico, e assim facilitar o agrupamento de espécies e evitar a criação de espécies muito raras geradas artificialmente. Utilizou-se um tamanho de janela de 5x5 para o filtro da média.

Figura 13. (a) Imagem original e seu histograma. (b) Imagem realçada e seu histograma. (c) Imagem suavizada e seu histograma.



Fonte: (ROCHA, 2014).

3.3 Extração de Características

Na fase de extração de características, inicialmente aplicamos a técnica de quantização, com o objetivo de agregar as informações de textura presentes em cada quantização e, assim, aumentar o poder discriminatório. As amostras realçadas foram quantizadas em 256, 128, 64, 32, 16 e 8 níveis de cinza. A partir de cada quantização, é calculado o LBP (*Local Binary Pattern*), para isso usamos como padrão janela 3 x 3. Sobre os LBPs extraímos os índices de diversidade de Gleason e Menhinick, para descrever a textura da amostra. Este cálculo é proposto através de quatro abordagens independentes: (1) a partir do histograma da imagem; (2) a partir da Matriz de Co-ocorrência de Níveis de Cinza (GLCM); (3) a partir da Matriz de Comprimentos de Corrida de Cinza (GLRLM); e (4) a partir da Matriz de Comprimentos de Lacuna de Cinza (GLGLM).

3.3.1 Índice de Diversidade Ecológica

Uma adaptação dos índices de diversidade ecológica é utilizada para gerar estatísticas para a análise de textura de imagens. Neste trabalho foram usados os índices de Gleason e Menhinick.

A adaptação do conceito de índice de diversidade foi feita a partir de duas abstrações. A primeira considera que uma comunidade será formada pelos LBPs da ROI, em que cada LBP é um indivíduo e o valor do LBP define a espécie. A segunda, na qual a comunidade é definida pelos elementos internos das matrizes de co-ocorrência de espécies e as espécies

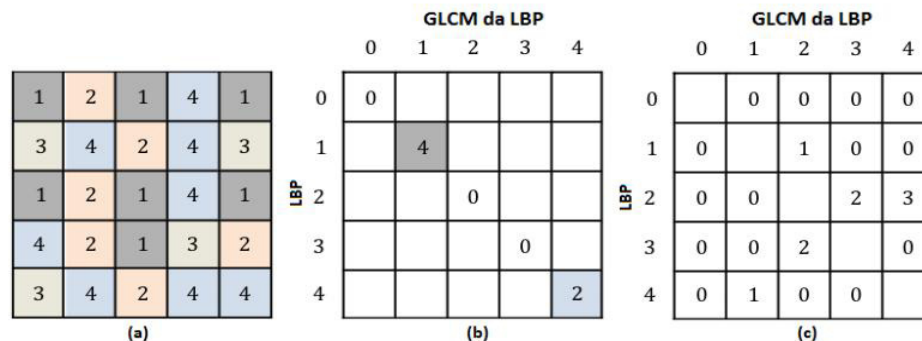
serão formadas pelos valores destes elementos. Assim, buscou-se verificar se nas ROIs existe a dominância de alguns padrões em relação a outros. O mesmo procedimento de extração de características foi aplicado aos dois índices.

A representação da ROI através do histograma possibilita a extração da riqueza de espécies (s) pela quantidade de entradas não nulas e a abundância relativa de cada espécie pelo valor dessas entradas, assim registrou-se a frequência de cada espécie. O vetor de características produzido apresenta 6 variáveis, calcula-se o valor da diversidade para cada quantização.

A matriz GLCM foi utilizada como forma de representação da ROI para verificar a diversidade da dominância de alguns pares de LBPs sobre outros. As comunidades de pares de LBPs são representadas de duas formas. Na primeira, os indivíduos correspondem às ocorrências de um par de LBPs (i, j) com o mesmo valor, separados por uma distância d e posicionados em uma direção θ . A diagonal principal da matriz GLCM representa a população de cada espécie. Na segunda, consideram-se como membros da comunidade as ocorrências de pares de LBPs (i, j) com valores diferentes. Assim, a população de indivíduos é representada fora da diagonal principal da matriz GLCM (Figura 14).

Para a direção θ foram adotados os valores $0^\circ, 45^\circ, 90^\circ$ e 135° . Para a distância dos valores utilizados foram 1, 2, 3, 4 e 5. O vetor de características apresentou 120 atributos de textura, 5 distâncias multiplicado por 4 direções multiplicado por 6 quantizações, e para cada θ e d é necessária uma GLCM e foram consideradas seis quantizações.

Figura 14. Cálculo da matriz GLCM para $\theta = 0^\circ$ e $d = 2$. (a) ROI 5 x 5. (b) Ocorrência de pares de LBPs de mesmo valor. (c) Ocorrência de pares de LBPs de valores diferentes.

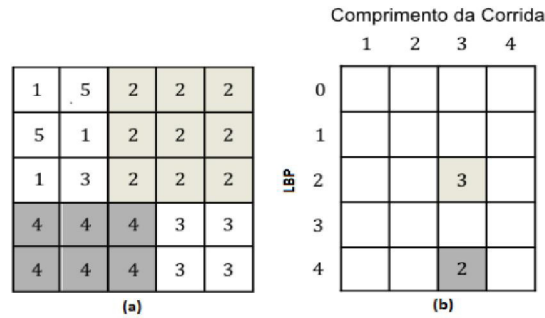


Fonte: (ROCHA, 2014)

A matriz GLRLM foi utilizada para analisar se há nas massas a predominância de corridas relativamente longas em relação às corridas curtas ou vice-versa. As comunidades

foram formadas pelas ocorrências de sequência consecutivas e colineares de n LBPs de mesmo valor em uma direção θ (Figura 15).

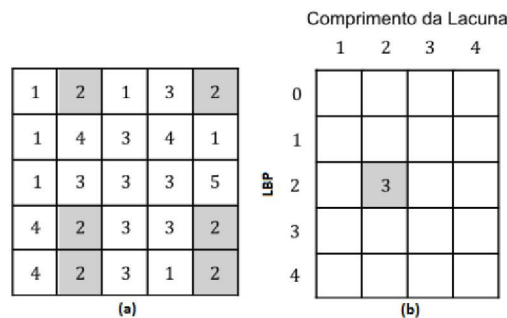
Figura 15. Cálculo da matriz GLRLM para $\theta = 0^\circ$. (a) ROI 5 x 5. (b) Ocorrência de corridas de LBPs de comprimento $k = 3$.



Fonte: (ROCHA, 2014).

A aplicação dos índices de diversidade com a matriz GLGLM visa investigar se uma massa apresenta, de uma maneira geral, a textura mais homogênea do que outra, é possível que ela contenha uma concentração maior de vizinhos homogêneos, sugerindo uma baixa diversidade. Se ocorrer o contrário, e apresentar uma menor concentração de vizinhos homogêneos, é provável que se tenha uma alta diversidade. As comunidades foram formadas por LBPs de valor i quando este LBP é encontrado apenas no início e no fim de uma sequência de LBPs consecutivos e colineares em uma direção θ (Figura 16).

Figura 16. Cálculo da matriz GLGLM para $\theta = 0^\circ$. (a) ROI 5 x 5. (b) Ocorrência de lacunas de LBPs de comprimento $k = 2$.



Fonte: (ROCHA, 2014).

Para a GLGLM e a GLRLM, foram consideradas seis quantizações, e quatro direções, 0° , 45° , 90° e 135° . Como é necessária uma matriz para cada quantização e direção, o vetor de características resultante apresenta 24 variáveis.

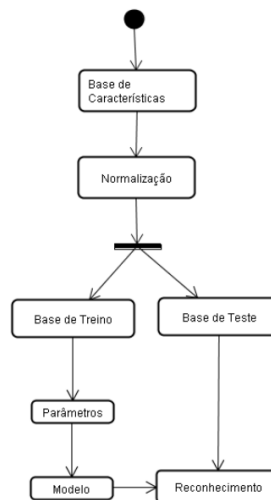
3.4 Reconhecimento de Padrões

A etapa final da metodologia proposta consiste em classificar as massas em maligna e benigna utilizando reconhecimento de padrões. As características de textura extraídas em etapas anteriores são submetidas aos classificadores supervisionados SVM e *Adaboost* para fazer tal classificação.

3.4.1 SVM (Support Vector Machine)

Durante a etapa de extração de características foram gerados vetores de características das amostras através do cálculo dos índices de diversidade de Gleason e Menhinick, a partir das quatro abordagens propostas (histograma, GLCM, GLRLM e GLGLM). De posse da base de características é necessário normalizar as diferentes características para uma faixa de valores comuns como -1 a 1. Esse mecanismo ajuda o classificador a convergir com maior facilidade na etapa de treinamento, e também padroniza a distribuição de valores das variáveis, as quais podem assumir diferentes domínios (BRAZ, 2008).

Figura 17. Fluxo de atividade da etapa de classificação.



Fonte: (BRAZ, 2008).

A base de características foi dividida randomicamente em dois grupos: base de treino e base de teste. Os percentuais usados neste trabalho para treino e teste foram respectivamente: 50/50, 60/40, 70/30 e 80/20. Para cada proporção foram realizadas 5 repetições do teste de forma aleatória. Como foi usado o núcleo radial do SVM, cada experimento teve os parâmetros de custo C e grau de complexidade da função de mapeamento γ . Os valores desses

parâmetros são estimados através de busca exaustiva realizada pelo script em *python grid.py*, pertencente ao pacote LIBSVM (CHANG, 2010). Este script busca, através de validação cruzada, a melhor combinação de parâmetros para a base, retornando o melhor percentual de acerto total sobre as amostras de treino e teste.

Durante a etapa de treinamento é gerado o modelo que o SVM utiliza para classificar as amostras de teste. O mecanismo de classificação, que desconhece as amostras de teste, busca se assemelhar com condições reais de teste, assim com o modelo gerado se torna possível realizar a etapa de reconhecimento de padrões com as amostras de teste separadas.

3.4.2 Adaboost

Para realizar os experimentos com o *Adaboost* utilizou-se a mesma base de características que foram aplicadas no SVM. O *Adaboost.M1* é uma adaptação do *Adaboost* original, implementado no pacote WEKA (*Waikato Environment for Knowledge Analysis*) e foi o algoritmo escolhido para ser utilizado neste trabalho. O pacote WEKA, desenvolvido pela Universidade de Waikato, na Nova Zelândia, é um software de domínio público que disponibiliza diversos algoritmos de aprendizado de máquina.

O *Adaboost* funciona em conjunto com outro algoritmo de aprendizado de máquina para fazer a classificação. Neste trabalho foi utilizado o SVM com o núcleo RBF, assim as características foram normalizadas para valores entre -1 e 1 e os parâmetros C e Y foram estimados usando o pacote LIBSVM. Da mesma forma como no experimento com o SVM, foram usadas 4 proporções para treinamento e teste da base de características, 50/50, 60/40, 70/30 e 80/20. Os testes foram realizados adotando 5 repetições para cada proporção e o parâmetro do número de iterações do *Adaboost.M1* foi definido em 10, valor padrão utilizado no WEKA.

3.5 Validação dos Resultados

A validação dos resultados é feita através do cálculo de algumas estatísticas sobre o resultado dos testes. Esta etapa permite medir o desempenho da metodologia como satisfatória ou não, sendo útil, inclusive, para apontar pontos positivos e negativos para melhoria futura deste trabalho. As estatísticas foram: Acurácia, Sensibilidade e Especificidade.

Essas métricas são comumente empregadas em sistemas CAD (*Computer-Aided Detection*) / CADx (*Computer-Aided Diagnosis*) e aceitas pela sociedade para a análise de desempenho de sistemas baseados em processamento de imagens.

4 RESULTADOS E DISCUSSÃO

Nesta seção serão apresentados os resultados dos testes obtidos com a metodologia proposta por este trabalho para a classificação de massas extraídas de imagens mamográficas em maligna e benigna. São demonstrados os resultados obtidos por cada índice de diversidade sobre o histograma e as matrizes GLCM, GLRLM e GLGLM.

4.1 Índices de Diversidade

Os resultados produzidos pelo experimento utilizando o índice de Gleason estão listados nas Tabelas 1 e 2. Para este índice utilizando o SVM, a abordagem que apresentou os melhores resultados foi a GLCM diagonal na proporção 50/50 com uma acurácia e sensibilidade de 77% e especificidade de 76%. Para o mesmo índice utilizando o Adaboost.M1 os melhores resultados também foram obtidos pela abordagem GLCM diagonal com a proporção de 60/40 e uma acurácia de 75, sensibilidade de 78% e especificidade de 71%.

Os resultados produzidos pelo experimento utilizando o índice de Menhinick estão listados nas Tabelas 3 e 4. Para este índice tanto o SVM quanto o Adaboost.M1 chegaram aos melhores resultados utilizando a abordagem GLCM Diagonal na proporção 50/50: o SVM apresentou uma acurácia de 77%, sensibilidade de 83% e especificidade de 71%; o Adaboost.M1 obteve uma acurácia de 75%, sensibilidade de 78% e especificidade de 72%.

Através da análise dos resultados foi possível constatar que os índices de diversidade apresentaram resultados muito semelhantes, o que pode ser explicado pelo fato de que ambos usam os mesmos parâmetros em seus cálculos. A análise também mostrou que a representação estatística da imagem usando a GLCM Diagonal foi a que apresentou os melhores resultados para os dois índices e para os dois classificados utilizados. Já o histograma apresentou os piores resultados no geral.

A Tabela 5 apresenta uma breve comparação entre os resultados encontrados neste trabalho e alguns trabalhos citados na Seção 1.3. A metodologia proposta apresentou resultados melhores em um caso e piores em outros, mas têm que se levar em conta as diferentes metodologias e número de amostras utilizadas nos experimentos.

Tabela 1: Resultados do Índice de Gleason utilizando SVM

Técnica	Proporção	Média Acurácia	Média Sensibilidade	Média Especificidade
GLCM Diagonal	50/50	77,50 ± 2,18	77,70 ± 1,51	76,55 ± 3,56
	60/40	76,15 ± 3,74	76,10 ± 1,59	75,45 ± 6,57
	70/30	75,88 ± 2,60	74,68 ± 2,78	77,12 ± 2,48
	80/20	76,45 ± 2,15	78,34 ± 3,21	75,96 ± 2,89
GLCM Matriz	50/50	75,58 ± 2,10	71,90 ± 1,90	80,01 ± 3,68
	60/40	74,34 ± 2,75	71,52 ± 2,26	75,21 ± 3,15
	70/30	75,78 ± 3,02	70,32 ± 2,75	80,60 ± 1,49
	80/20	75,61 ± 2,36	70,93 ± 4,45	80,11 ± 5,91
GLGLM	50/50	69,15 ± 2,13	67,10 ± 3,63	70,29 ± 3,10
	60/40	68,31 ± 2,65	68,47 ± 4,85	66,71 ± 7,90
	70/30	70,01 ± 1,97	70,92 ± 4,55	71,20 ± 3,15
	80/20	67,90 ± 4,50	68,18 ± 2,30	68,87 ± 4,01
GLRLM	50/50	76,51 ± 1,31	75,48 ± 2,12	76,10 ± 0,98
	60/40	72,01 ± 1,27	70,29 ± 1,91	74,02 ± 1,28
	70/30	76,97 ± 2,78	73,25 ± 3,16	79,93 ± 2,03
	80/20	73,45 ± 3,39	68,69 ± 3,34	77,78 ± 3,68
Histograma	50/50	70,13 ± 2,08	70,65 ± 2,10	70,09 ± 6,71
	60/40	70,47 ± 1,02	68,59 ± 1,90	70,88 ± 3,01
	70/30	71,50 ± 1,69	71,14 ± 3,01	72,31 ± 2,86
	80/20	72,06 ± 3,17	72,90 ± 3,66	72,59 ± 3,71

Tabela 2: Resultados do Índice de Gleason utilizando Adaboost.M1

Técnica	Proporção	Média Acurácia	Média Sensibilidade	Média Especificidade
GLCM Diagonal	50/50	75,15 ± 0,87	74,96 ± 0,83	75,46 ± 1,31
	60/40	75,62 ± 0,92	78,78 ± 1,78	71,98 ± 2,50
	70/30	73,18 ± 2,46	75,30 ± 3,45	70,82 ± 3,00
	80/20	74,62 ± 2,03	79,50 ± 4,03	69,34 ± 5,28
GLCM Matriz	50/50	71,88 ± 1,26	70,58 ± 2,55	73,48 ± 2,27
	60/40	70,16 ± 1,95	70,06 ± 3,80	70,22 ± 2,18
	70/30	70,30 ± 2,89	71,10 ± 2,89	69,40 ± 3,59
	80/20	74,62 ± 2,03	79,50 ± 4,03	69,34 ± 5,28
GLGLM	50/50	67,90 ± 0,82	66,90 ± 2,97	69,12 ± 2,64
	60/40	68,76 ± 1,92	66,74 ± 3,61	71,28 ± 0,52
	70/30	69,18 ± 2,84	69,86 ± 2,31	68,48 ± 4,81
	80/20	68,08 ± 2,54	68,50 ± 4,38	67,76 ± 1,89
GLRLM	50/50	71,82 ± 1,73	74,84 ± 2,71	68,28 ± 2,05
	60/40	72,08 ± 1,36	74,48 ± 2,16	69,34 ± 2,97
	70/30	73,28 ± 2,57	75,02 ± 3,49	71,18 ± 4,66
	80/20	72,32 ± 5,39	74,96 ± 4,04	69,18 ± 7,15
Histograma	50/50	67,34 ± 1,78	68,58 ± 2,41	65,96 ± 3,79
	60/40	68,50 ± 1,75	72,48 ± 2,12	63,76 ± 1,37
	70/30	65,70 ± 1,11	68,06 ± 2,45	63,26 ± 4,71
	80/20	65,98 ± 2,47	67,70 ± 3,77	64,56 ± 8,00

Tabela 3: Resultados do Índice de Menhinick utilizando SVM

Técnica	Proporção	Média Acurácia	Média Sensibilidade	Média Especificidade
GLCM Diagonal	50/50	77,50 ± 4,96	83,93 ± 2,15	71,36 ± 7,77
	60/40	74,04 ± 1,83	80,09 ± 2,79	68,41 ± 5,20
	70/30	76,92 ± 1,54	80,81 ± 2,67	70,97 ± 3,68
	80/20	75,77 ± 1,17	79,46 ± 3,04	72,30 ± 4,31
GLCM Matriz	50/50	73,28 ± 2,95	74,35 ± 4,69	71,98 ± 9,26
	60/40	75,15 ± 1,51	74,62 ± 2,76	73,11 ± 3,71
	70/30	73,43 ± 1,89	74,73 ± 5,78	72,79 ± 6,35
	80/20	73,98 ± 1,37	73,41 ± 2,53	73,63 ± 2,03
GLGLM	50/50	69,35 ± 0,81	71,98 ± 1,95	66,14 ± 1,59
	60/40	68,70 ± 1,29	71,75 ± 3,34	65,33 ± 4,00
	70/30	70,39 ± 1,01	74,13 ± 3,88	68,62 ± 2,46
	80/20	69,18 ± 4,09	72,63 ± 4,46	66,81 ± 4,93
GLRLM	50/50	75,47 ± 0,94	77,19 ± 1,09	73,43 ± 1,10
	60/40	76,95 ± 1,90	78,90 ± 1,70	74,78 ± 4,84
	70/30	77,06 ± 1,73	78,85 ± 2,52	74,73 ± 1,39
	80/20	77,14 ± 2,19	77,39 ± 3,41	77,06 ± 4,81
Histograma	50/50	70,62 ± 0,86	75,80 ± 1,07	61,13 ± 1,03
	60/40	70,99 ± 0,71	72,42 ± 4,89	67,56 ± 4,94
	70/30	68,93 ± 2,59	74,68 ± 2,55	63,12 ± 6,87
	80/20	73,34 ± 1,02	73,49 ± 3,01	73,25 ± 4,66

Tabela 4: Resultados do Índice de Menhinick utilizando Adaboost.M1

Técnica	Proporção	Média Acurácia	Média Sensibilidade	Média Especificidade
GLCM Diagonal	50/50	75,46 ± 1,98	78,04 ± 1,95	72,40 ± 3,01
	60/40	75,04 ± 1,48	77,82 ± 2,17	71,50 ± 4,68
	70/30	74,74 ± 2,21	79,12 ± 5,29	69,98 ± 3,30
	80/20	73,34 ± 2,21	77,66 ± 4,29	68,60 ± 5,62
GLCM Matriz	50/50	72,70 ± 0,75	74,18 ± 2,66	70,98 ± 3,60
	60/40	71,10 ± 1,69	71,46 ± 4,17	70,62 ± 5,13
	70/30	68,90 ± 1,36	69,88 ± 2,84	67,68 ± 3,37
	80/20	72,06 ± 3,61	72,14 ± 3,26	71,68 ± 6,59
GLGLM	50/50	67,90 ± 2,33	65,72 ± 2,65	70,52 ± 2,69
	60/40	68,78 ± 0,85	68,14 ± 3,97	69,86 ± 5,25
	70/30	67,98 ± 1,44	67,20 ± 4,05	69,14 ± 3,96
	80/20	67,08 ± 1,54	68,08 ± 2,01	66,06 ± 2,79
GLRLM	50/50	75,36 ± 1,29	77,24 ± 2,20	73,16 ± 2,67
	60/40	74,16 ± 1,49	76,26 ± 1,15	71,66 ± 2,39
	70/30	75,54 ± 1,20	76,62 ± 3,50	74,26 ± 2,73
	80/20	74,72 ± 2,56	76,58 ± 1,67	72,90 ± 4,91
Histograma	50/50	68,48 ± 0,56	70,72 ± 4,25	65,98 ± 3,86
	60/40	68,92 ± 0,83	73,10 ± 2,27	64,06 ± 0,90
	70/30	68,02 ± 1,54	71,52 ± 2,54	64,04 ± 1,66
	80/20	68,56 ± 3,68	72,76 ± 3,48	64,42 ± 6,57

Tabela 5: Comparação com alguns trabalhos referentes à classificação de massas em imagens de mamografias em maligno e benigno.

Trabalhos	Base de Dados	Acurácia
(ROCHA, 2014)	DDSM	92%
(LIU, 2011)	DDSM	66%
(NANNI, 2012)	DDSM	88%
Nossa Metodologia	DDSM	77%

5 CONCLUSÃO

Neste trabalho apresentamos a viabilidade da utilização da combinação de abordagens estrutural e estatística para a análise de textura, e do SVM e Adaboost.M1 para classificação em maligna e benigna de massas presentes em mamografia.

A metodologia testou uma abordagem estrutural, o *Local Binary Pattern* (LBP) e 2 abordagens para extração de características, os índices de diversidade de Gleason e Menhinick, combinados com a representação da imagem através de estatísticas de primeira ordem (histograma), segunda ordem (GLCM) e ordem superior (GLGLM e GLRLM).

Os índices de diversidade apresentaram resultados semelhantes o que pode ser explicado pelo fato de os dois levarem em conta em seus cálculos os mesmos parâmetros. O índice de diversidade de Gleason obteve seu melhor resultado utilizando a técnica GLCM Diagonal nos dois classificadores, no SVM a proporção 50/50 apresentou uma acurácia de 77%, no Adaboost.M1 uma acurácia de 76% para a proporção 60/40. O índice de diversidade de Menhinick também obteve seu melhor resultado utilizando a técnica GLCM Diagonal, a proporção 50/50 com acurácia de 77% também foi a melhor para os dois classificadores.

O histograma dentre todas as formas de representação estatística da imagem utilizada neste trabalho, foi a que obteve os piores resultados.

Em relação aos classificadores utilizados, SVM e Adaboost.M1, o SVM na maioria dos testes apresentou resultados melhores que o Adaboost.M1.

Os resultados obtidos evidenciaram que a análise de textura como parâmetro para a diferenciação entre os padrões malignidade e benignidade de massa em imagens de mamografias é importante, entretanto uma tarefa difícil, principalmente devido ao fato de que é muito comum massas malignas e benignas possuírem características de textura semelhantes.

Como propostas futuras para encaminhamento deste trabalho, podemos citar o uso de máscaras maiores e formas modificadas do LBP, para tentar analisar de uma forma diferente a estrutura da textura, também propomos fazer uso de outros índices de diversidade. Outra forma de complementar este trabalho é a utilização de outros modelos *Ensemble* para a classificação e também testar o Adaboost com outras técnicas de aprendizagem de máquina como classificador parcial.

REFERÊNCIAS BIBLIOGRÁFICAS

ACS, A. C. S. Learn About Breast Cancer. Disponível em: <http://www.cancer.org>. Último Acesso: 13/10/2014. 2014.

BEBIS, G. Advances in Visual Computing. *Lecture Notes in Computer Science (LNCS)*, LNCS, v. 4291, 4292, 2006.

BLAND, M. *An introduction to medical statistics*. New York: Oxford University Press, 2000.

BRAZ, JR., G. *Classificação de Regiões de Mamografias em Massa e Não Massa usando Estatística Espacial e Máquina de Vetores de Suporte*. Dissertação (Mestrado) – Curso de Pós-Graduação em Engenharia de Eletricidade da Universidade Federal do Maranhão, São Luís – MA, 2008.

BROWER, J. E.; ZAR, J. H.; ENDER, C. V. *Field and Laboratory Methods for General Ecology*. Dubuque: Mcgraw-hill College, 1997.

CARVALHO, P. M. S. *Classificação de Tecidos de Mama a Partir de Imagens Mamográficas em Massas e Não Massas Usando Índice de Diversidade de Mcintosh e Máquina de Vetores de Suporte*. Dissertação (Mestrado) - Curso de Pós-Graduação em Engenharia de Eletricidade da Universidade Federal do Maranhão, São Luís – MA, 2012.

CHANG, C.; LIN, C. LIBSVM – A Library for Support Vector Machines. *ACM Transactions on Intelligent Systems and Technology*, v. 2, n. 3, p. 27-27, 2011. Disponível em: <http://www.csie.ntu.edu.tw/~cjlin/libsvm/>.

CHAVES, A. C. F. *Extração de Regras Fuzzy para Máquinas de Vetor de Suporte (SVM) para Classificação em Múltiplas Classes*. Tese (Doutorado) - Pontifícia Universidade Católica do Rio de Janeiro, Rio de Janeiro, 2006.

DUARTE, J. C. *O Algoritmo Boosting at Star e Suas Aplicações*. Tese (Doutorado) - Pontifícia Universidade Católica do Rio de Janeiro, Rio de Janeiro, 2009.

FREUND, Y.; SCHAPIRE, R. E. A decision-theoretic generalization of on-line learning and an application to boosting. In: *European Conference on Computational Learning Theory*. [s.l.: s.n.], 1995. p. 23–37

ROCHA, S. V. *Diferenciação do Padrão de Malignidade e Benignidade de Massas em Imagens de Mamografias Usando Padrões Locais Binários, Geoestatística e Índice de Diversidade*. Tese (Doutorado) - Curso de Pós-Graduação em Engenharia de Eletricidade da Universidade Federal do Maranhão, São Luis - 2014.

GALLOWAY, M. M. Texture Analysis using Gray Level Run Lengths. *Computer Graphics and Image Processing*, v. 4, p. 172-179, 1975.

GONZALEZ, R.; WOODS, R. *Digital Image Processing*. Addison-Wesley Reading, Mass, 1992.

GONZALEZ, R.; WOODS, R. *Digital image processing*. New Jersey: Pearson Prentice Hall, 2002.

GONZALEZ, R.; WOODS, R. *Processamento Digital de Imagens*. 3a.ed. São Paulo: Pearson Prentice Hall, 2010.

HARALICK, R. M.; SHANMUGAM, K.; DINSTEN, I. Textural Features for Image Classification. *IEEE Transactions on Systems, Man, and Cybernetics*, v. 3, n. 6, p. 610-621, 1973.

HAYKIN, S.; ENGEL, P. *Redes Neurais: Princípios e Prática*. Porto Alegre: Bookman, 2001.

INCA. Instituto Nacional do Câncer. Disponível na internet em: <http://www.inca.gov.br/estimativa/2014/>. Acesso em: 16 jul. 2014.

LIU, X.; LIU, J.; TANG, J. Improved local binary patterns for classification of masses using mammography. *Systems, Man, and Cybernetics (SMC), 2011 IEEE International Conference on*, IEEE, p. 2692-2695, 2011.

LOONEY, C. *Pattern Recognition using Neural Networks: Theory and Algorithms for*

Engineers and Scientists. New York: Oxford University Press, Inc, 1997.

MAGURRAN, A. E. *Measuring Biological Diversity*. Padstow, U.K: Blackwell Science, 2004. 248 p.

MARQUES, O.; VIEIRA, H. *Processamento Digital de Imagens*. Rio de Janeiro: Brasport, 1999.

MASCARO, A. A.; MELLO, C. A. B.; P., S. W.; CAVALCANTI, G. D. C. Mammographic images segmentation using texture descriptors. *31st Annual International Conference of the IEEE EMBS*, IEEE, p. 3653-3656, 2009.

MENHINICK, E. F. A comparison of some species-individuals diversity indices applied to samples of field insects. *Ecology*, v. 45, n. 4, p. 859-861, 1964.

MORSE, B. *Data Structures for Image Analysis*. 2000. Brigham Young University. Disponível em: http://homepages.inf.ed.ac.uk/rbf/CVonline/LOCAL_COPIES/MORSE/data-structures.pdf. Último Acesso em: 02/09/2014

NASCIMENTO, L. B. *Classificação de Nódulos Pulmonares em Maligno e Benigno utilizando os Índices de Diversidade de Shannon e de Simpson*. Dissertação (Mestrado em Engenharia de Eletricidade). Universidade Federal do Maranhão 2012.

NANNI, L.; BRAHNAM, S.; LUMINI, A. A very high performing system to discriminate tissues in mammograms as benign and malignant. *Expert Systems with Applications*, Elsevier, v. 39, n. 4, p. 1968-1971, 2012.

NANNI, L.; LUMINI, A.; BRAHNAM, S. Local Binary Pattern variantes as texture descriptors for medical image analysis. *Artificial Intelligence in Medicine*, Elsevier, v. 49, n. 4, p. 117-125, 2010.

OJALA, T.; PIETIKAINEN, M.; HARWOOD, D. A. A comparative study of texture measures with classification based on feature distribution. *Pattern Recognition*, v.29, n. 1, p 51-59, 1996.

PEDRINI, H.; SCHWARTZ, W. R. *Análise de Imagens Digitais: Princípios, Algoritmos e Aplicações*. [S.l.]: Thomson Learning, 2008. 503 p.

SANTOS, V. K. *Uma generalização da distribuição do índice de diversidade generalizada por Good com aplicação em Ciências Agrárias*. 57 p. Dissertação (Mestrado) – Universidade Federal Rural de Pernambuco – UFPE. Recife, 2009

VAPNIK, V. N. *Statistical Learning Theory*. New York: Wiley, 1998.

XINLI, W.; ALBREGTSEN, F.; FOYN, B. Texture Features from Gray Level Gap Length Matrix. *Workshop on Machine Vision Applications*, p. 375-378, 1994.