

DEEP LEARNING PARA DIAGNÓSTICO ANATOMOPATOLÓGICO DE CÂNCER

Anotações
14-09-2022

Código atual: github.com/dan-teix/classificador_cancer

Resumo: Classifica subtipos de câncer pulmonar do *LC25000 dataset* utilizando a CNN de [arXiv:2009.03878](https://arxiv.org/abs/2009.03878)

Problemas: Abordagem insuficiente para problemas na medicina

1. IA na medicina

- **Diagnóstico:** algoritmo identifica e determina a doença, subtipo, localização de tumores etc.
- **Prognóstico:** prevê evolução da saúde do paciente, e.g. risco de morte baseado em exames laboratoriais e outros dados do paciente.
- **Tratamento:** modelos que estimem o efeito de determinado tratamento num paciente.

2. Principais problemas e possíveis soluções

- **Classes desbalanceadas:** muito mais imagens de pacientes saudáveis criam a tendência no algoritmo de estimar uma probabilidade pequena para presença de doença.
 - Soluções: modificar a função de perda para atribuir diferentes pesos para as diversas classes — maior peso para as classes menos representadas; reamostragem para balancear as classes.
- **Tamanho dos conjuntos de dados:** número limitado de exemplos comparativamente com classificação de outros tipos de imagens.
 - Soluções: *data augmentation*, e.g. ruídos nas cores em histopatologia; utilizar redes neurais pré-treinadas em outros conjuntos de imagens e aplicar aprendizagem por transferência — primeiras camadas de uma rede aprendem aspectos mais genéricos, que podem ser comuns a outros problemas, enquanto a últimas camadas aprendem características mais específicas (nessas inclui-se *fine-tuning*).
- **Exames de um mesmo paciente:** modelo pode “memorizar” uma imagem utilizada na fase de treino e usar essa informação em imagens futuras na fase de teste.
 - Soluções: incluir todos os exames de um mesmo paciente no mesmo conjunto, seja de treino ou validação.
- **Amostragem:** em conjuntos de dados pequenos, ao extrair o conjunto de testes do mesmo, pode ser que o último não contenha amostras da classe minoritária.
 - Soluções: garantir que o conjunto de teste tenha pelo menos alguma porcentagem (e.g. 50%) da classe minoritária e manter mesma distribuição no conjunto de validação. Caso isso resulte num conjunto de treino desbalanceado, ver “classes desbalanceadas” acima.
- **Desacordo entre avaliadores:** para gerar rótulos com padrão de referência, pode-se estabelecer um consenso através de votos majoritários por especialistas; exames adicionais que possam levar a um diagnóstico definitivo.

3. Métricas

Sensibilidade (taxa de verdadeiros positivos) e **especificidade** (taxa de verdadeiros negativos) são mais adequadas do que acurácia. Tais métricas seguem da **matriz de confusão**. Incluir curva ROC e AUC.

4. Além da classificação: segmentação de imagem

Traçar quais partes da imagem contém uma doença — importante na quantificação do tamanho de tecidos, na localização de tumores e no planejamento de tratamentos.

Estratégia para representar imagens 3D de ressonância magnética (MRI): combinar as múltiplas sequências tratando-as como diferentes canais (análogo a *rgb*, mas não limitado a três canais).

Possível problema: sequências de MRI desalinhadas — requer **registro das imagens**.

Segmentação: identificar as bordas de diversos tecidos (ou de tumores), ou seja, determinar a classe de cada ponto (voxel) no volume 3D das sequências de MRI combinadas:

- Abordagem 2D: folhear o volume 3D com fatias 2D que, por sua vez, vão passar por um modelo de segmentação (arquitetura U-Net 2D) sendo, ao final, novamente combinadas para fornecer a segmentação do volume inicial.
 - Ponto fraco: perda de informação devido ao folheamento.
- Abordagem 3D: particionar o volume inicial em subvolumes 3D suportados pela memória ao passar pelo algoritmo de segmentação (arquitetura U-Net 3D).
 - Ponto fraco: assim como no caso anterior, há perda de informações topológicas.

⚠ *Data augmentation*: aplicar mesmas operações aos outputs da segmentação.

Função de perda: *soft dice loss* é popular em algoritmos de segmentação.

5. Dificuldades de implementação de técnicas de IA

- Diferenças populacionais e tecnológicas
- Validação externa
- *Black box*