

# Entregável 1.2. Diferenças entre os resultados obtidos por PyCCD e eeCCD

## Índice

1	Introdução .....	2
2	Dados .....	3
3	CCD.....	4
3.1	Adaptação do algoritmo PyCCD .....	4
3.2	Parâmetros do PyCCD .....	5
3.3	Escolha do valor do parâmetro de regularização <i>alpha</i> .....	6
4	Desempenho do CCD sobre dados de referência .....	7
5	Referências .....	9

## Índice das Figuras

Figura 1– Disponibilidade de dados Theia/MAJA e GEE/s2cloudless. ....	3
Figura 2 – Localização de uma amostra de 10.000 pixels para input do PyCDD. ....	4
Figura 3 – Resultado obtido pelo CCD para um pixel onde houve uma ocorrência de fogo na data identificada.....	5
Figura 4 – Desempenho do modelo com diferentes valores de alpha para dados Theia/MAJA (gráficos superiores) e GEE/s2cloudless (gráficos inferiores). ....	7
Figura 5 – Classificação de detecções automáticas pelo PyCCD em relação às identificações dos analistas. VP (Verdadeiro Positivo), VN (Verdadeiro Negativo), FP (Falso Positivo), FN (Falso Negativo). O retângulo rosa representa a tolerância na incerteza temporal (30 dias) – cf. (Moraes et al., ND).....	8

## Índice das Tabelas

Tabela 1 – Parâmetros escolhidos para o CCD.....	5
Tabela 2 – Erros de omissão e comissão para dados Theia/MAJA e GEE/s2cloudless. ....	8
Tabela 3 – Desempenho dos dois algoritmos: PyCCD (usando Theia/MAJA e GEE/s2cloudless) e eeCCD (usando apenas GEE/s2cloudless). ....	9

# 1 Introdução

O desenvolvimento de metodologias eficazes à escala nacional e eficientes a nível computacional para a criação sistemática de um produto nacional em formato vetorial de delimitação de manchas superiores a 0.5 ha de perda recente de floresta e mato com base em análise automática de imagens de satélite requer ferramentas que permitem explorar a informação temporal e espacial contida nos dados de muito grande dimensão que consistem em séries de imagens Sentinel-2 para Portugal Continental.

As técnicas que parecem ser mais adaptadas ao problema incluem algoritmos de deteção de alterações ao nível do pixel, entre os quais se destaca o designado *Continuous Change Detection and Classification* (CCDC) que não só permite identificar perdas de vegetação (Tarefa 3) como caracteriza o padrão temporal da ocupação do solo com o objetivo de caracterizar o fator de alteração (Tarefa 4).

Neste relatório, descreve-se as adaptações realizadas e os resultados de experiências computacionais realizadas com o CCDC sobre uma base de dados de referência para:

1. Adaptação do algoritmo PyCCD (aberto) a dados Sentinel-2;
2. Parametrização do método;
3. Avaliação do desempenho;
4. Comparação entre o desempenho do PyCCD e o eeCCD.

Os resultados do presente relatório permitem:

1. Aplicar CCDC para deteção automática de alterações a uma nova coleção de dados de referência (dados Navigator), e posterior validação;
2. Definir direções de pesquisa no quadro do projeto para permitir reduzir o tempo de execução do método com o objetivo de o aplicar à generalidade do território continental português, em áreas onde se prevê à partida a possibilidade de ocorrência de perdas de vegetação e mato.

Recorde-se que CCDC é uma técnica de deteção e classificação de mudanças contínuas que tem sido amplamente utilizada em monitorização ambiental, especialmente na deteção de mudanças na cobertura terrestre ao longo do tempo (Zhu et al., 2014). Esta técnica é particularmente útil em aplicações de deteção remota, onde é necessário identificar e quantificar mudanças em grandes áreas de forma contínua e precisa (Awty-Carroll et al., 2019).

O CCDC é projetado para lidar com séries temporais de imagens de deteção remota, como imagens de satélite, que são adquiridas em intervalos regulares ao longo do tempo. A ideia fundamental do CCDC é modelar as séries temporais de dados de deteção remota de forma a capturar tanto a variabilidade sazonal quanto as mudanças de longo prazo na cobertura terrestre. Isso é feito através da combinação de técnicas de modelos de estatística, como séries temporais de decomposição e modelos de regressão (Xian et al., 2022, Zhu et al., 2014).

Os scripts em Python, objeto deste relatório, estão disponíveis em: [https://github.com/manuelcampagnolo/S2CHANGE/tree/main/scripts/pyccd\\_theia/notes](https://github.com/manuelcampagnolo/S2CHANGE/tree/main/scripts/pyccd_theia/notes) (como já mencionado no Entregável E.1.1).

## 2 Dados

Os testes realizados e descritos neste relatório usaram imagens Sentinel-2 fornecidos pela Theia, uma plataforma francesa que desempenha um papel importante na distribuição e processamento de dados de detecção remota. Os dados disponibilizados pela Theia são pré-processados e corrigidos usando o algoritmo de correção atmosférica MAJA, proporcionando uma base sólida para análises subsequentes. Além disso, foram realizados testes com dados do Sentinel-2 provenientes do *Google Earth Engine* (GEE), os quais foram processados utilizando o algoritmo de correção atmosférica s2cloudless (Skakun et al., 2022).

Ambos os conjuntos de dados foram criteriosamente selecionados, descartando imagens com cobertura de nuvens superior a 60%, visando garantir a qualidade e a consistência dos resultados obtidos.

Na Figura 1, apresenta-se a série temporal de disponibilidade de imagens para os dados Theia/MAJA (representados por círculos vermelhos) e os dados GEE/s2cloudless (representados por círculos azuis) entre 2017 e o final de 2021.

Observa-se que no início da série, os dados Theia/MAJA são mais consistentes em comparação com os GEE/s2cloudless. No entanto, para datas mais recentes os dados GEE/s2cloudless têm maior densidade relativamente aos dados Theia/MAJA, o que pode ser explicado por serem ambas máscaras conservativas, mas MAJA usar um buffer de maior dimensão para nuvens (240 m em vez de 160 m) e incluir detecção de sombras (Skakun et al., 2022).

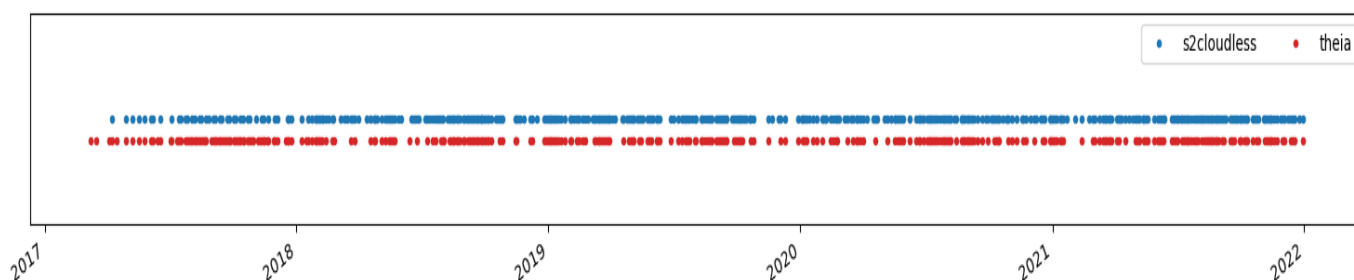


Figura 1– Disponibilidade de dados Theia/MAJA e GEE/s2cloudless.

Estes testes foram realizados sobre a base de dados de referência BDR-300 (Moraes et al., ND) extraída sobre o tile T29TNE Sentinel-2 para o período compreendido entre março de 2018 e dezembro de 2021. Esta região de Portugal Continental é caracterizada pela presença de densas florestas e vastas áreas agrícolas, representando um ambiente diversificado. Além disso, esta área foi severamente afetada por incêndios florestais de grande magnitude em 2017, resultando numa paisagem em constante mudança. Assim, dada a complexidade e a dinâmica desse cenário, o PyCCD demonstra um forte potencial para detetar e monitorizar mudanças nesta região.

Para a maiorias das experiências realizadas e descritas neste relatório, foram escolhidos aleatoriamente 10000 pixels entre os aproximadamente 350000 pixels existentes em BDR-300. A Figura 2 mostra a distribuição espacial de uma mostra de 10000 pixels extraídos da BDR-300, sendo que o número dentro de cada hexágono representa a quantidade de pixels escolhidos dentro do respetivo buffer. A distribuição dos pontos foi feita de maneira a assegurar que todos os buffers tivessem pelo menos um ponto para

serem incorporados nos dados de entrada do algoritmo de detecção de alterações (Moraes et al., ND).

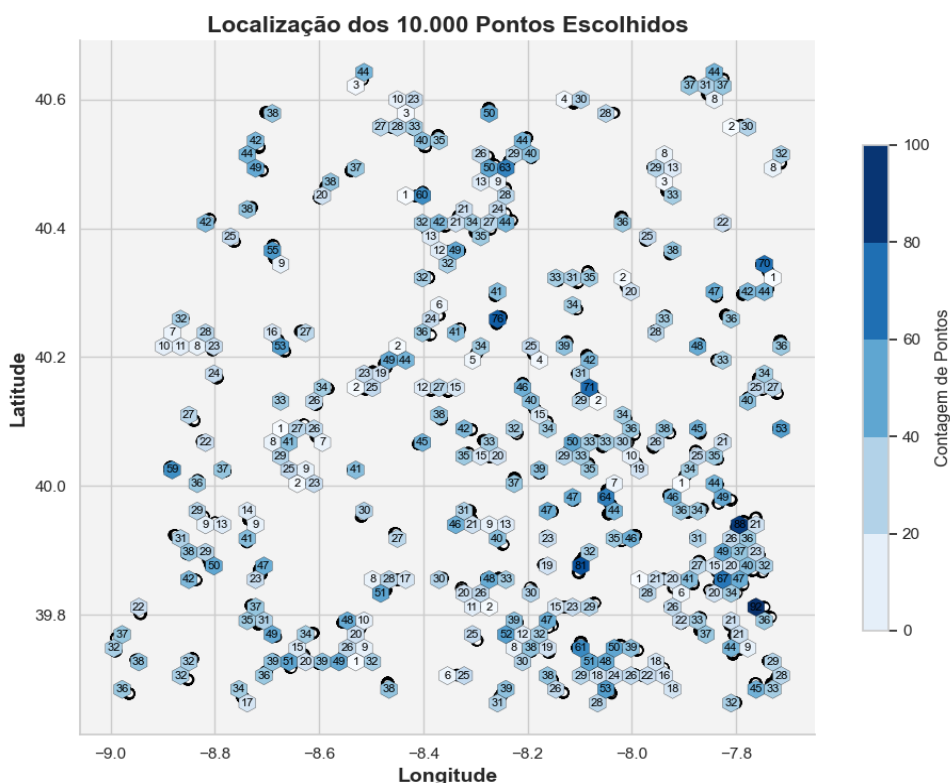


Figura 2 – Localização de uma amostra de 10.000 pixels para input do PyCDD.

## 3 CCD

### 3.1 Adaptação do algoritmo PyCCD

Para implementar eficientemente o CCDC, foi feita uma adaptação de um código Python disponível em <https://code.usgs.gov/lcmap/pyccd>, originalmente projetado para dados do satélite LANDSAT, mas que teve de ser adaptado, no quadro deste projeto, para suportar dados do satélite Sentinel-2. Esse algoritmo é designado por PyCCD.

O código do PyCCD aberto é muito complexo e não está documentado em detalhe, sendo a interpretação do código e a sua melhoria uma tarefa tecnicamente difícil. No entanto, entende-se que esse investimento é necessário para poder dispor para a prossecução do projeto de um código que possa correr em diversos ambientes computacionais, desde workstations locais a clusters de computação avançada como, por exemplo, os clusters da Rede Nacional de Computação Avançada (RNCA). A escolha de um código na linguagem Python facilita a portabilidade do código entre plataformas.

Durante o processo de adaptação do código, ajustamos o modelo para incorporar as diferentes características dos dois satélites. Por exemplo, o LANDSAT possui uma banda térmica, ao passo que o Sentinel-2 não. Além disso, os dados do Sentinel-2 utilizados no algoritmo já incluíam correções atmosféricas, diferentemente dos dados do LANDSAT, onde esse processamento era feito dentro do próprio algoritmo.

Adicionalmente, foram adicionados alguns parâmetros ao modelo que anteriormente não estavam incorporados, como '**CHISQUAREPROB**' e '**ALPHA**' (este parâmetro determina a intensidade da regularização aplicada ao modelo).

### 3.2 Parâmetros do PyCCD

Os parâmetros principais escolhidos para o modelo PyCCD foram os seguintes:

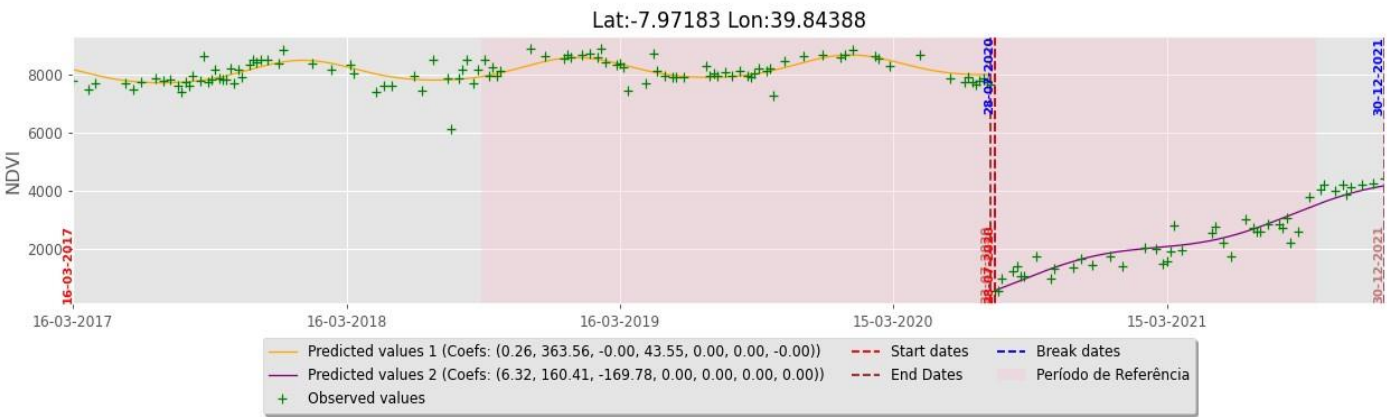
PEEK_SIZE	MIN_YEARS	DETECTION_BANDS	LASSO_MAX_ITER	CHISQUAREPROB	ALPHA
6	1	NDVI, Green, SWIR2	1000, 25000	0.999	0 a 200

*Tabela 1 – Parâmetros escolhidos para o CCD.*

Os valores dos parâmetros foram selecionados com base nos resultados de (Moraes et al., ND) mas adaptados para os novos dados de entrada, em particular as imagens Theia. Nesse estudo, os autores testaram o CCD implementado em GEE com os mesmos valores para os parâmetros mencionados anteriormente, atingindo um bom desempenho do modelo.

Foram aqui considerados 3 combinações de dados e de algoritmos: imagens Theia/MAJA com PyCCD, imagens GEE/s2cloudness com algoritmo PyCCD, e imagens GEE/s2cloudness com algoritmo eeCCD, o algoritmo que é executado em GEE e cujo código não é aberto. O objetivo foi verificar que as soluções locais (usando PyCCD) tinham o mesmo nível de desempenho do que o produto de alterações obtido em ambiente GEE, com o algoritmo eeCCD.

A Figura 3 ilustra o resultado obtido quando o CCD é implementado num determinado pixel. Inicialmente, os dados revelam uma série temporal consistente de valores de NDVI, variando entre 7500 e 8500, indicando uma cobertura vegetal saudável e densa, predominantemente composta por pinheiro-bravo. Entretanto, a 28 de julho de 2020 (designada por data de quebra) observa-se uma mudança abrupta. Os valores de NDVI sofrem uma queda acentuada, atingindo a faixa de 1000. Essa redução significativa e repentina nos valores do NDVI é um claro indicador de uma situação adversa, especificamente um evento de fogo.



*Figura 3 – Resultado obtido pelo CCD para um pixel onde houve uma ocorrência de fogo na data identificada.*

### 3.3 Escolha do valor do parâmetro de regularização *alpha*

Para melhorar o desempenho do modelo, foi feito uma série de testes para determinar o valor ideal do parâmetro *alpha*, que desempenha um papel fundamental na função do Lasso da biblioteca Python scikit-learn, determinando a intensidade da regularização aplicada ao modelo. Quando este parâmetro é definido como zero, a regressão Lasso é desativada, resultando numa regressão linear.

Esta pesquisa foi realizada em ambos os conjuntos de dados, Theia/MAJA (apresentado nos gráficos superiores da Figura 4) e GEE/s2cloudless (exibido nos gráficos inferiores da Figura 4), explorando uma ampla faixa de valores de *alpha*, variando de 0 a 200.

Ao investigar os diferentes conjuntos de dados, observamos que o valor ideal de *alpha* pode variar para cada um dos conjuntos de dados, por várias razões (cf. documentação da função `sklearn.linear_model.lasso`):

- I. Tamanho do conjunto de dados: Conjuntos de dados de diferentes tamanhos podem exigir diferentes níveis de regularização para se ajustar aos padrões específicos de cada conjunto.
- II. Distribuição dos dados: Variações na distribuição dos dados entre os conjuntos podem afetar o comportamento do modelo, levando a necessidades distintas de regularização.
- III. Complexidade dos dados: Conjuntos de dados mais complexos, ou com maior presença de ruído, podem requerer valores de *alpha* diferentes para controlar adequadamente o ajuste do modelo.

Além disso, a escolha do valor de *alpha* influencia o equilíbrio entre sobreajustamento e subajustamento do modelo. Valores mais altos de *alpha* tendem a reduzir o sobreajustamento, resultando em modelos mais simples com menos coeficientes significativos. Por outro lado, valores mais baixos de *alpha* oferecem maior flexibilidade ao

modelo, mas podem aumentar o risco de sobreajustamento, levando a modelos mais complexos com mais coeficientes significativos.

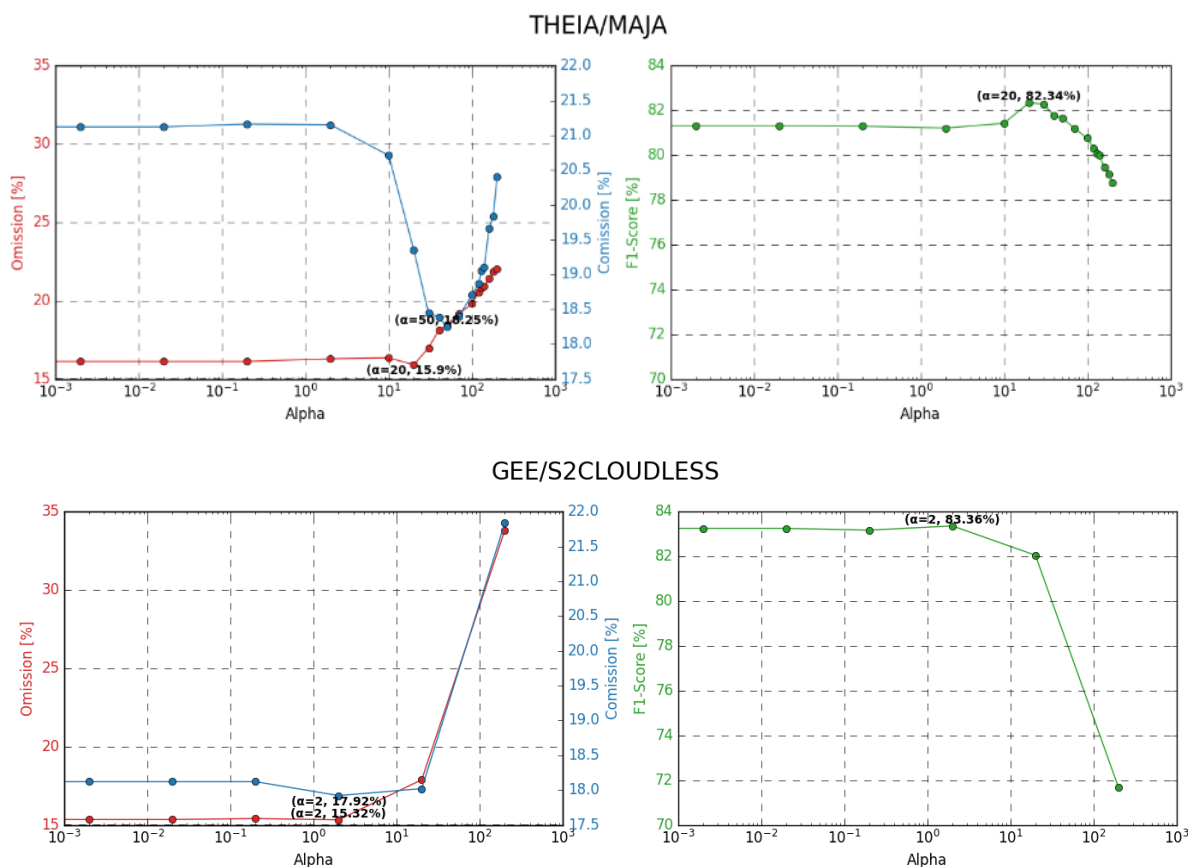


Figura 4 – Desempenho do modelo com diferentes valores de  $\alpha$  para dados Theia/MAJA (gráficos superiores) e GEE/s2cloudless (gráficos inferiores).

Ao analisar os dados Theia/MAJA, o valor ideal de  $\alpha$  para maximizar o F1-Score e minimizar a Omissão é de 20, resultando em taxas de 82.34% e 15.9%, respetivamente. Embora o valor que minimiza o erro de comissão seja 50, a escolha recaiu sobre 20 devido à mínima discrepância observada nessa métrica entre os valores de 20 e 50, resultando numa taxa de 18.38% (em vez dos 18.25% que seriam obtidos com  $\alpha = 50$ ).

Em relação aos dados GEE/s2cloudless, o valor ideal de  $\alpha$  para otimizar tanto o F1-Score quanto a Omissão e a Comissão é 2, produzindo taxas de 83.36%, 15.32% e 17.92%, respetivamente.

## 4 Desempenho do CCD sobre dados de referência

Após o modelo ter processado uma amostra de pixels, é realizado um teste de validação utilizando os resultados fornecidos pelos analistas. Esse teste de validação abrange o período de análise compreendido entre 12/09/2017 e 30/09/2021, no qual cada detecção automática feita pelo PyCCD foi classificada como verdadeiro positivo (VP), verdadeiro negativo (VN), falso positivo (FP) ou falso negativo (FN) – Figura 5.

- **Verdadeiro Positivo (VP):** Refere-se a uma detecção automática pelo PyCCD que coincide exatamente com a quebra identificada pelo analista.

- **Verdadeiro Negativo (VN):** Pode ocorrer de duas maneiras:
  1. A detecção pelo PyCCD ocorre fora do período de análise, portanto, durante este período, não há detecções identificadas nem pelo PyCCD nem pelos analistas.
  2. O PyCCD não identifica nenhuma detecção ao longo da série temporal.
- **Falso Positivo (FP):** Ocorre quando o PyCCD detecta algo automaticamente, mas o analista não identifica nenhuma quebra.
- **Falso Negativo (FN):** Pode surgir de duas maneiras:
  1. O PyCCD não identifica nenhuma detecção automática ao longo de toda a série, enquanto o analista identifica.
  2. O PyCCD identifica antes do período de análise dos analistas.

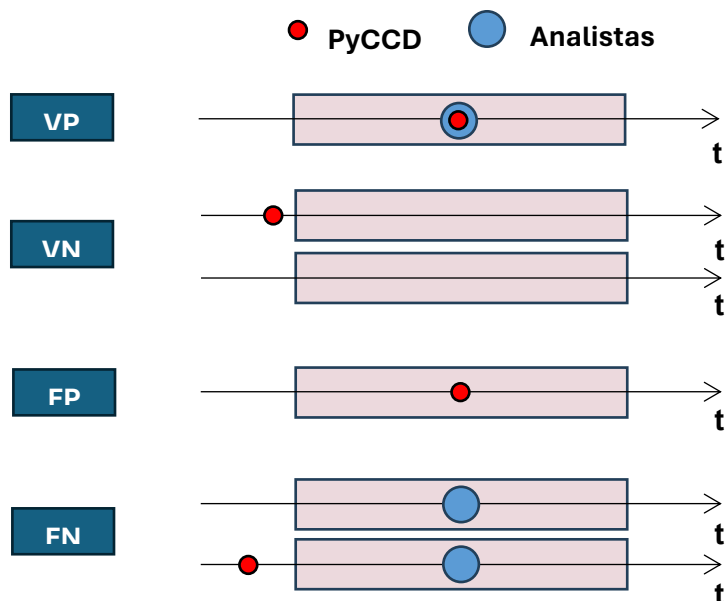


Figura 5 – Classificação de detecções automáticas pelo PyCCD em relação às identificações dos analistas. VP (Verdadeiro Positivo), VN (Verdadeiro Negativo), FP (Falso Positivo), FN (Falso Negativo). O retângulo rosa representa a tolerância na incerteza temporal (30 dias) – cf. (Moraes et al., ND).

Relacionando estas definições com os erros de omissão e comissão calculados para os dados Theia/MAJA e GEE/s2cloudless, podemos entender melhor a precisão das detecções efetuadas pelo CCD em relação às identificações realizadas pelos analistas (Tabela 2).

	THEIA/MAJA	GEE/S2CLOUDLESS
<b>Erro de omissão [%]</b>	15.90	15.32
<b>Erro de comissão [%]</b>	18.38	17.92

Tabela 2 – Erros de omissão e comissão para dados Theia/MAJA e GEE/s2cloudless.

Os resultados dos erros de omissão para os dados Theia e GEE revelam que o algoritmo PyCCD enfrenta dificuldades na detecção precisa de mudanças em ambas as bases de dados, perdendo cerca de 15.9% e 15.32% das mudanças reais identificadas pelos analistas, respectivamente. Da mesma forma, os erros de comissão indicam que aproximadamente 18.38% das mudanças identificadas pelo PyCCD nos dados Theia/MAJA e cerca de 17.92% nos dados GEE/s2cloudless não correspondem a mudanças reais observadas pelos analistas. Estes resultados terão de ser comparados com os resultados sobre a nova base de dados de referência a ser explorada no seguimento do projeto.



Refira-se que a comparação dos resultados GEE/s2cloudless e GEE/eeCCD (com o algoritmo de código não aberto do GEE) deram resultados muito semelhantes (Tabela 3), o que é uma indicação de o que o desempenho de PyCCD é comparável ao que se obtém usando o GEE, tendo a vantagem de poder ser executado em plataformas locais.

	<b>F1-Score [%]</b>	<b>Omission Error [%]</b>	<b>Commission error [%]</b>
<b>PyCCD</b>			
<b>THEIA/MAJA</b>	82.34	15.90	18.38
<b>GEE/s2cloudless</b>	83.36	15.32	17.92
<b>eeCCD</b>			
<b>GEE/s2cloudless</b>	81.14	15.77	21.73

*Tabela 3 – Desempenho dos dois algoritmos: PyCCD (usando Theia/MAJA e GEE/s2cloudless) e eeCCD (usando apenas GEE/s2cloudless).*

## 5 Referências

1. Awty-Carroll, Katie; Bunting, Pete; Hardy, Andy; Bell, Gemma (2019). Using Continuous Change Detection and Classification of Landsat Data to Investigate Long-Term Mangrove Dynamics in the Sundarbans Region. Remote Sensing, 11(23), 2833–. doi:10.3390/rs11232833
2. Código PyCCD para dados LANDSAT: <https://code.usgs.gov/lcmap/pyccd> (Acedido janeiro de 2024)
3. Moraes, D. et al., "Continuous vegetation loss monitoring in a dynamic landscape of Central Portugal with Sentinel-2 data", (ND) em processo de revisão para o International Journal of Applied Earth Observation and Geoinformation
4. Skakun, Sergii & Wevers, Jan & Brockmann, Carsten & Doxani, Georgia & Aleksandrov, Matej & Batic, Matej & Frantz, David & Gascon, Ferran & Gómez-Chova, Luis & Hagolle, Olivier & López-Puigdollers, Dan & Louis, Jerome & Lubej, Matic & Mateo-Garcia, Gonzalo & Osman, Julien & Peressutti, Devis & Pflug, Bringfried & Puc, Jernej & Richter, Rudolf & Žust, Lojze. (2022). Cloud Mask Intercomparison eXercise (CMIX): An evaluation of cloud masking algorithms for Landsat 8 and Sentinel-2. Remote Sensing of Environment. 274. 112990. 10.1016/j.rse.2022.112990.
5. Sklearn.linear\_model.lasso (no date) Scikit-learn: Machine Learning in Python. Disponível em: [https://scikit-learn.org/stable/modules/generated/sklearn.linear\\_model.Lasso.html](https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.Lasso.html) (Acedido a: 29 de abril de 2024)
6. Xian, G. Z., Smith, K., Wellington, D., Horton, J., Zhou, Q., Li, C., Auch, R., Brown, J. F., Zhu, Z., & Reker, R. R. (2022). Implementation of the CCDC algorithm to produce the LCMAP Collection 1.0 annual land surface change product. Earth System Science Data, 14(1), 143–162. <https://doi.org/10.5194/essd-14-143-2022>
7. Zhu, Zhe; Woodcock, Curtis E. (2014). Continuous change detection and classification of land cover using all available Landsat data. Remote Sensing of Environment, 144(), 152–171. doi:10.1016/j.rse.2014.01.011