

# **Tomografia PCA aplicado às galáxias do Projeto CALIFA Survey**

Eduardo Alberto Duarte Lacerda

Orientador:

Prof. Dr. Roberto Cid Fernandes Jr.

...

Universidade Federal de Santa Catarina  
Centro de Ciências Físicas e Matemáticas  
Curso de Pós-Graduação em Física

Dissertação de mestrado apresentada ao Curso de Pós-Graduação em Física da UFSC em preenchimento parcial dos requisitos para obtenção do título de Mestre em Física.

Florianópolis (SC) – 8 de fevereiro de 2014

# Sumário

<b>1</b>	<b>Introdução</b>	<b>1</b>
1.1	Um labirinto de dados . . . . .	1
1.2	A nova ferramenta - Tomografia PCA . . . . .	2
1.3	Este trabalho . . . . .	3
1.3.1	Organização deste trabalho . . . . .	4
<b>2</b>	<b>O projeto CALIFA e o <i>pipeline</i> PyCASSO</b>	<b>6</b>
2.1	O survey CALIFA . . . . .	6
2.1.1	A “ <i>colméia</i> ” de fibras - os dados do CALIFA usados neste trabalho . . . . .	7
2.2	O <i>pipeline</i> PyCASSO . . . . .	9
2.2.1	E/S de dados no PyCASSO . . . . .	9
2.2.2	Exemplos de utilização . . . . .	10
<b>3</b>	<b>PCA e Tomografia PCA</b>	<b>14</b>
3.1	Principal Component Analysis . . . . .	14
3.1.1	PCA das galáxias do CALIFA . . . . .	15
3.2	Tomografia PCA . . . . .	17
3.2.1	Evidências de linhas largas . . . . .	17
<b>4</b>	<b>PCA nos espectros</b>	<b>21</b>

---

4.1	Pré-processamento dos cubos . . . . .	21
4.1.1	Normalização . . . . .	23
4.1.2	Cinemática . . . . .	24
4.1.3	Linhas de emissão e intervalos específicos em comprimento de onda .	24
4.1.4	Fluxos observados e sintéticos . . . . .	27
<b>5</b>	<b>Resultados</b>	<b>29</b>
<b>6</b>	<b>Conclusões e perspectivas</b>	<b>30</b>
6.1	Este trabalho . . . . .	30
6.2	Trabalhos futuros . . . . .	30
<b>A</b>	<b>[FIXME] Python</b>	<b>31</b>
	<b>Referências Bibliográficas</b>	<b>I</b>

# Lista de Figuras

2.1	Diagrama cor-magnitude para as galáxias do CALIFA. . . . .	8
2.2	Configuração do <i>bundle</i> de fibras do PPMAS/PPAK. . . . .	8
2.3	Exemplo de programa utilizando PyCASSO . . . . .	11
2.4	Programa idade estelar média . . . . .	11
2.5	Mapa da idade estelar média da galáxia NGC 2916 . . . . .	12
2.6	Exemplo de programa para perfil radial . . . . .	12
2.7	Perfil radial da idade estelar média da galáxia NGC 2916 . . . . .	13
3.1	Exemplo de cálculo de PCA usando o PyCASSO e SciPy . . . . .	16
3.2	<i>Scree test</i> na galáxia NGC 4736. . . . .	18
3.3	Tomograma e autoespectro 1 da galáxia NGC 4736. . . . .	19
3.4	Tomograma e autoespectro 2 da galáxia NGC 4736. . . . .	19
3.5	Tomograma e autoespectro 3 da galáxia NGC 4736. . . . .	20
4.1	Exemplo de máscaras em um espectro do cubo de dados. . . . .	22
4.2	Fluxos de normalização para cada zona da galáxia K0277. . . . .	23
4.3	Tomogramas de 1 a 4 da galáxia NGC2916 - sem normalização. . . . .	25
4.4	Tomogramas de 1 a 4 da galáxia NGC2916 - com normalização. . . . .	26

# Lista de Tabelas

2.1	Relação de pixels e zonas em algumas galáxias do CALIFA . . . . .	10
-----	---	----

# Capítulo 1

## Introdução

### 1.1 Um labirinto de dados

Cientistas hoje, em sua maioria, encontram-se perdidos em meio a um labirinto de informações. Essas são apenas parte de um cenário incompleto, mas não falso, que chamamos de nosso universo. Com o avanço tecnológico, melhores formas de se obter informações tornam cada vez mais evidente o surgimento de um *fordismo*<sup>1</sup> relativo a informações que auxiliam na criação desse labirinto, mas ao mesmo tempo fomenta a criatividade e curiosidade de cientistas que se tornam aventureiros na busca da saída. Mas essa fuga é apenas a primeira etapa assim se tornando parte basilar na formação da pesquisa científica.

Os primeiros *surveys*<sup>2</sup> astronômicos surgem com a inata curiosidade do homem de observar tudo a sua volta e registrar suas observações. O crescimento desses catálogos é produto direto da evolução dos equipamentos usados nestas investigações. No início da década de 80 (Huchra et al. 1983; Huchra 1988; da Costa et al. 1988) a astronomia extragalática entra nesse cenário de produção sistemática em massa de dados. De lá para cá a quantidade de dados só aumenta e, com a criação dos *mega-surveys* (*SDSS*; York et al. 2000) (*2dFGRS*; Colless 1999) (*2MASS*; Skrutskie et al. 2006), alguns já terminados, outros ainda começando ou por terminar (*LSST*; Ivezić et al. 2008) (*J-PAS*; Benítez et al. 2009), estamos à beira do humanamente impossível de se avaliar, necessitando assim a ajuda de máquinas e métodos computacionais cada vez mais eficientes.

---

<sup>1</sup>O termo *fordismo* foi criado por Antonio Gramsci em 1922 e é relativo a Henry Ford e o surgimento da produção em massa de automóveis no início do século XX.

<sup>2</sup>Um *survey* astronômico é um levantamento de informações ou mapeamento de regiões do céu utilizando telescópios e detetores.

Com esse crescimento exponencial na quantidade de dados, precisamos cada vez mais de ferramentas matemáticas/estatísticas. No contexto de espectros de galáxias, uma ferramenta muito útil para síntese espectral é o STARLIGHT, desenvolvido por Cid Fernandes et al. (2005). Hoje, com o uso de painéis de fibras óticas apontadas para as galáxias temos os *surveys* de IFS (*Integral Field Spectroscopy*) onde passamos a obter dezenas, centenas e até milhares de espectros por galáxia, obtendo então dimensão espacial também nos dados. Assim temos para cada píxel (duas dimensões espaciais) um espectro (uma dimensão espectral), formando assim um cubo de *spaxels*<sup>3</sup>. O pioneiro nessa produção em massa de dados é o *Calar Alto Legacy Integral Field spectroscopy Area survey*<sup>4</sup> (CALIFA; Sánchez et al. 2012), produzindo cerca de mil espectros por galáxia observada. Outros *mega-surveys* IFS estão por vir (veja a seção 2.1).

Com esses cubos de dados em mão, podemos assim executar o STARLIGHT para cada *spaxel* e então obtemos propriedades físicas em função da posição na galáxia (Cid Fernandes et al. 2013). Uma ferramenta menos astrofísica, mas não menos utilizada, é o PCA (*Principal Component Analysis*) e no presente trabalho fazemos seu uso juntamente com uma técnica criada por Steiner et al. (2009) chamada de Tomografia PCA, onde une-se imageamento e espectrografia ao mesmo tempo.

## 1.2 A nova ferramenta - Tomografia PCA

A técnica de Análise de Componentes Principais (PCA) é simples, não-paramétrica e nos ajuda a extrair informações de conjuntos de dados com muitas variáveis, reduzindo a dimensionalidade no sentido de encontrar quais são os elementos com maior variâncias no seu conjunto de dados. No caso de um espectro, temos uma medida de fluxo para cada intervalo em comprimento de onda. Em um cubo de *spaxels*, temos as dimensões espaciais também, obtendo assim uma infinidade de variáveis. O grande problema do PCA é que você tem a resposta, mas não sabe a pergunta.

Hoje, PCA é utilizado exaustivamente em várias áreas de conhecimento, principalmente em reconhecimento de padrões, computação visual e filtragem e compactação de dados (Kamruzzaman et al. 2010; Borcea et al. 2012). Podemos ver exemplos também em medicina (Balakrishnan et al. 2013). Na astrofísica o PCA passeia por diversas ramificações da área. A luz de um objeto até nossos telescópios sofre influência de muito ruído, devido a inúmeros

---

<sup>3</sup>spectral pixels

<sup>4</sup><http://www.caha.es/CALIFA/>

problemas como: atenuações, avermelhamento por poeira, contaminação ótica através de objetos que estejam no mesmo FoV, entre outros; a diversa quantidade de instrumentos e suas complexidades também geram diversas assinaturas indesejadas; por todos esses motivos fica claro que os dados geralmente necessitam de uma boa filtragem. Junto com outras técnicas (tomogramas, wavelets, Fourier), o PCA vem sendo muito utilizado para filtragem de dados, principalmente cubos de dados advindos de IFS que possuem muitas dessas assinaturas instrumentais (Riffel et al. 2011). Outro exemplo de uso de PCA aparece no artigo IV do grupo SEAGal/STARLIGHT (Mateus et al. 2007) auxiliando no estudo da dependência ambiental de algumas propriedades físicas (idade estelar média ponderada pela luz, massa estelar, metalicidade estelar e *mass-to-light ratio* ( $M_{\star}/L$ )) em uma amostra de galáxias do SDSS DR4. Em Chen et al. (2012) é criada uma biblioteca com 25 mil modelos de espectros de galáxias com diferentes idades, metalicidades, velocidade de dispersão, SFH (*star formation history*), extinção por poeira e aplicado PCA em cima dessa biblioteca. Usando uma minimização quadrática encontram quais os coeficientes e qual o número ótimo de PCs que melhor estimam os parâmetros físicos dos espectros modelo. Então projetam os espectros observados pelo SDSS DR7 (Abazajian et al. 2009) e pelo Baryon Oscillation Spectroscopic Survey (BOSS Ahn et al. 2012) para galáxias massivas com  $z \sim 0.6$  até o presente, atribuindo um sentido físico a cada PC.

Com a criação da técnica de Tomografia PCA citada anteriormente temos, além da análise das componentes principais, uma imagem formada pela matriz de covariância (as questões matemáticas serão abordadas no Capítulo 3). Assim podemos saber o peso, ou relevância, de cada componente principal (daqui pra frente PC) relacionada a uma posição na galáxia, aliando a infinidade de informações físicas presentes nos espectros e nas imagens. Além do artigo de Steiner et al. (2009), outros exemplos de uso de Tomografia PCA podem ser encontrados em Riffel et al. (2011) e Ricci et al. (2011).

## 1.3 Este trabalho

Através da colaboração do Grupo de Astrofísica da Universidade Federal de Santa Catarina (GAS-UFSC) com o grupo de pesquisadores do projeto CALIFA temos a oportunidade de trabalhar com os dados de IFS das galáxias observadas por esse projeto, que ainda está em andamento. O seu primeiro *Data Release*<sup>5</sup> (Husemann et al. 2013, DR1) possui 100 objetos e por volta de 400 mil espectros. A previsão é que ao término do projeto serão observadas

<sup>5</sup>[http://www.caha.es/CALIFA/public\\_html/?q=content/califa-dr1](http://www.caha.es/CALIFA/public_html/?q=content/califa-dr1)



até 600 objetos. Embora outros surveys de IFS estão completos ou em andamento (ver Seção 2.1), o CALIFA é o que podemos chamar de *estado da arte* em surveys de espectroscopia de campo integrado (IFS).

Dado o grande número de informações sobre cada galáxia necessitamos de um *pipeline* que faça a organização de todos os dados para que cálculos e gráficos de mais variadas dificuldade sejam fácil para que até um programador de nível iniciante possa fazer. André L. de Amorim, colaborador de nosso grupo, juntamente com outros colaboradores de nosso grupo e do projeto CALIFA construiu o PyCASSO (*Python CALIFA STARLIGHT Synthesis Organizer*) (Cid Fernandes et al. 2013, cap. 4) que faz a organização dos dados que vêm do *survey* juntamente com a síntese de populações estelares feitas com o STARLIGHT, facilitando, e muito, o trabalho de quem usa estes dados. Sem a ajuda deste organizador, este trabalho seria muito mais difícil e acredito que não seria realizável em tempo hábil.

### 1.3.1 Organização deste trabalho

No seguinte capítulo apresenta de forma mais detalhada o *survey* CALIFA e o *pipeline* PyCASSO. Nele também são demonstrados alguns exemplos de utilização e gráficos.

O terceiro capítulo descreve matematicamente a técnica PCA e a Tomografia PCA, bem como sua utilização no presente trabalho.

A partir do quarto capítulo temos os mais diferentes *tipos*<sup>6</sup> de execução do PCA, juntamente com suas implicações aos resultados das análises. Também são apresentados diversos gráficos demonstrando as diferenças entre a utilização de cada pré-processamento, análise da Tomografia PCA e das PCs, bem como as principais referências na área.

Com todo o arcabouço teórico em mãos, no capítulo cinco temos um estudo de caso para 9 objetos do survey CALIFA, escolhidos por seu tipo morfológico. São 2 galáxias *early-type*, 4 espirais e 2 objetos compostos (mergers). Como forma de comparação com os resultados da síntese de populações estelares executadas pelo STARLIGHT nesses objetos também é feito uma espécie de engenharia reversa através de correlações e comparações.

Por fim, temos as conclusões e perspectivas futuras deste trabalho e da nossa colaboração com o projeto CALIFA no sexto e último capítulo. [!oj!] Devido ao grande número de

---


<sup>6</sup> [!oj!] na literatura usa-se *flavors*, não sei se devo usar *sabores* também. tiposs aqui é usado apenas como uma maneira de tipificar diferentes pré-processamentos dos dados antes da execução da análise de componentes principais.

imagens e dados para análise, muitas delas ficaram em anexo.

## Capítulo 2

# O projeto CALIFA e o *pipeline* PyCASSO

As observações do universo modificaram completamente o nosso modo de viver, pensar, compreender-se. Aprendemos a contar os dias, desenvolvemos um sistema de meses, estações do ano, movimentos das marés, entre outras coisas que já são parte do senso comum, mas que um dia foram o estado da arte da ciência. É assim que surge o CALIFA: um projeto que está modificando nossa maneira de ver e pensar as galáxias no nosso universo de forma que entendamos melhor a nossa também.



Com a massiva quantidade de dados obtidos, resultado direto de um projeto de ciência de ponta, vem também a dificuldade da interpretação dos dados. No caso do CALIFA, através da *pipeline* PyCASSO (Cid Fernandes et al. 2013), a programação investigativa se torna simples e ao mesmo tempo robusta, facilitando a construção de todo tipo de resultados, sejam eles  matemáticos ou físicos.

### 2.1 O survey CALIFA

No sul da Espanha, mais precisamente em *Sierra de Los Filabres* (Andalucía), está situado o germano-espânico *Calar Alto Observatory*. O projeto CALIFA está sendo possível através de observações pelo maior de seus 3 telescópios (3.5m) ao longo de 250 noites. Em comparação com o *SDSS*, o CALIFA terá a mesma ordem de número de espectros para estudo ( $\sim 10^6$ ), mas, apesar de um número menor de galáxias, graças ao IFU será o com melhor completeza por objeto. Existem alguns poucos surveys IFU e todos com, além de poucos objetos e FoV

menor, focos de estudo muito estreitos, dificultando o legado do survey para outras pesquisas científicas mais abrangentes (SAURON; de Zeeuw et al. 2002, região central de 72 galáxias com  $z < 0.01$ .) (PINGS; Rosales-Ortega et al. 2010, algumas galáxias muito próximas ( $\sim 10$  Mpc) e o estudo atual de 70 (U)LIRGs com  $z < 0.26$ ) (VENGA; Blanc et al. 2010, 30 galáxias espirais). Apesar de ser primariamente construído para o estudo da física bariônica da evolução de galáxias, o CALIFA está projetado para que seu legado seja bem abrangente, possibilitando diversos tipos de estudos em diversas áreas. Outros surveys IFU ainda estão por vir, como SAMI (Croom et al. 2012) e MaNGA<sup>1</sup>

### 2.1.1 A “colméia” de fibras - os dados do CALIFA usados neste trabalho

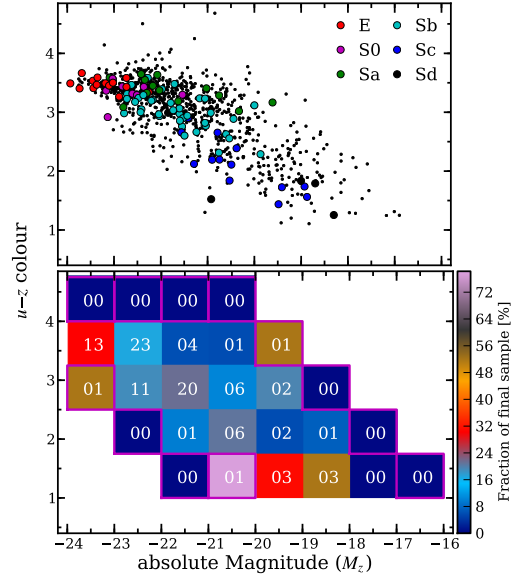
A amostra-mãe do projeto comporta 939 galáxias (das quais  $\sim 600$  serão observadas) com *redshifts* entre  $0.005 < z < 0.03$  que distribuídas cobrem o diagrama cor-magnitude com  $M_r < -18$  (Figura 2.1) em uma ampla variedade de tipos morfológicos, massa em estrelas, condições do gás ionizante. Para melhor aproveitar o *FoV* (*field-of-view*<sup>2</sup>) do instrumento de IFU é feito também um corte em dimensão ( $\sim 1'$  em diâmetro). No telescópio usado para o projeto está instalado o equipamento Potsdam Multi Aperture Spectrograph (PMAS; Roth et al. 2005) no modo PPAK (Verheijen et al. 2004; Kelz et al. 2006) formando um espectrofotômetro de campo integrado com um *bundle* de 382 fibras (Figura 2.2), das quais, 331 são para observação dos objetos, outras 36 para  *sky background sample* e outras 15 para calibração. As *science fibers* (331) cobrem um campo de visão hexagonal de  $74'' \times 64''$  que, através de uma técnica de três pontos de dithering  torna possível a observação de 100% do campo.

Os dados são reduzidos utilizando o programa CALIFA Pipeline versão 1.3c, descrito em Husemann et al. (2013). Os espectros vêm em duas configurações: a V500 cobrindo de  $\sim 3700$  até  $7000 \text{ \AA}$  com resolução de  $\sim 6 \text{ \AA}$  de largura à meia altura (FWHM) e a V1200 ( $\sim 3650 - 4600 \text{ \AA}$  FWHM  $\sim 2.3 \text{ \AA}$ ). A cobertura do V500 seria ideal para os propósitos de ciência feita pelo STARLIGHT mas por problemas com *vignetting* com a parte azul dessa configuração, os dados são reamostrados numa combinação das duas, criando uma que chamamos de COMBO. A parte com  $\lambda < 4600 \text{ \AA}$  vem do V1200 e a outra parte do V500. Os espectros foram reamostrados no mesmo FWHM do V500.

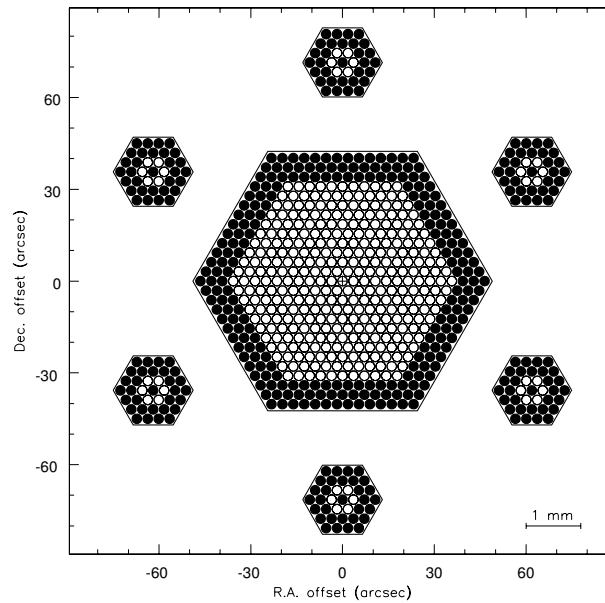
Após os dados passarem pelo CALIFA Pipeline v1.3c, os parâmetros físicos aqui usados

<sup>1</sup><http://www.sdss3.org/future/manga.php>

<sup>2</sup>campo de visão





**Figura 2.1:** Distribuição das galáxias do CALIFA no diagrama  $u - z$  vs.  $M_z$ . *Painel superior:* Em pontos pretos estão as galáxias pertencente a amostra-mãe e em cores as galáxias presente no CALIFA DR1. As diferentes cores representam os diferentes tipos morfológicos. *Painel inferior:* A fração de galáxias observadas pelo DR1 em relação a amostra-mãe. Retirado de Husemann et al. (2013), figura 2.



**Figura 2.2:** Este é o esquema com o *bundle* hexagonal com as 331 fibras de observação e mais 36 de amostra de céu. Retirado de Verheijen et al. (2004), figura 5.

vêm do PyCASSO que organiza as saídas da síntese de populações estelares resultantes da execução de cada espectro do IFU pelo STARLIGHT (Cid Fernandes et al. 2005) como descrito em Cid Fernandes et al. (2013) e também na próxima seção.

## 2.2 O pipeline PyCASSO

A beleza da espectroscopia de campo integral é poder unir imageamento e espectroscopia, obtendo-se assim a galáxia vista como um cubo de dados  $(x, y, \lambda)$ ; para cada  $\lambda$  temos uma imagem, e para cada par de ascensão reta e declinação  $(x, y)$  temos um espectro.  [Pedir uma imagem chique de cubo para o André](#). Apesar de cada espectro poder ser analisado individualmente, nos cubos do CALIFA existe uma correlação entre pixels vizinhos devido a  *seeing* do céu e ao processo de observação. Como nas regiões mais afastadas do núcleo da galáxia o brilho superficial é menor, há uma queda na relação sinal-ruído (S/N) dos dados. De maneira que haja S/N constante, é feito um agrupamento de pixels (zonificação de Voronoi) nas regiões mais afetadas (ver (Cid Fernandes et al. 2013, sec. 3). Vale aqui frisar que a grande maioria das zonas comporta um pixel apenas. Para a galáxia NGC 2916, por exemplo, 93% das zonas possuem apenas 1 pixel e  $\sim 3\%$  com mais de 10 pixels agrupados, mesmo com a imposição de relação S/N em 20 (veja Tabela 2.1). Portanto o cubo original é transformado numa matriz de zonas e comprimentos de onda. Com esses dados é então possível realizar a síntese de populações estelares com o STARLIGHT, resolvendo espacialmente as propriedades físicas estelares das galáxias (Cid Fernandes et al. 2013, 2014; González Delgado et al. 2013).

### 2.2.1 E/S de dados no PyCASSO

O PyCASSO é uma biblioteca desenvolvida em Python para organizar os dados da síntese feita pelo STARLIGHT. A versão usada neste trabalho é a 0.9.3<sup>3</sup>. Para um acesso mais rápido e reutilizável do código e dos dados em qualquer ambiente, organiza os cubos em formatos FITS ou HDF5. Em outra camada, várias matrizes e cubos são armazenados para acesso com nomes próprio (um exemplo, `popx`, representa a fração de luz distribuída pelas populações estelares) de forma que a programação exploratória não precise se preocupar com as características de cada formato de armazenamento de dados. Por fim, existe uma camada construída para análise, com funções que retornam a indexação de cada zona para um par  $(x, y)$ , cálculos de

<sup>3</sup><http://minerva.ufsc.br/~andre/PyCASSO-0.9.3/>

**Tabela 2.1:** Relação entre números de pixels por zona nas galáxias do CALIFA utilizadas neste trabalho.  $N_z$  representa o número de zonas.  $N_1$  é o numero de zonas com 1 pixel apenas e  $N_{10}$  aquelas que possuem mais de 10 pixels por zona.

Nome da galáxia	CALIFA ID	Hubble Type	$N_z$	$N_1$	$N_1/N_z$	$N_{10}$	$N_{10}/N_z$
NGC 0001	K0008	Sbc	1132	1077	0.95	40	0.04
NGC 0776	K0073	Sb	1733	1628	0.94	61	0.04
NGC 1167	K0119	S0	1879	1771	0.94	50	0.03
NGC 2623	K0213	Scd	561	530	0.94	19	0.03
NGC 2916	K0277	Sbc	1638	1528	0.93	53	0.03
NGC 4210	K0518	Sb	1938	1847	0.95	38	0.02
ARP 220	K0802	Sd	1157	1103	0.95	39	0.03
NGC 6515	K0864	E3	887	811	0.91	44	0.05

perfis radiais e azimutais, geometria, entre outras rotinas. Um programador facilmente pode adicionar mais rotinas como essas.

### 2.2.2 Exemplos de utilização

Ler um arquivo FITS é fácil com o PyCASSO, assim como o acesso aos dados. Na Figura 2.3 temos um exemplo de leitura de arquivo FITS e um cálculo da idade estelar média da galáxia a partir da idade média ponderada pela luminosidade, por zona (`at_flux__z`). A idade estelar média é calculada usando a expressão  $\langle \log t \rangle_L^{gal} = \sum_z \langle \log t \rangle_{L,z} L_z / \sum_z L_z$  onde  $L_z$  é a luminosidade e  $\langle \log t \rangle_L$  é a idade estelar média, ambas por zona.

Usando o `matplotlib`<sup>4</sup> podemos fazer gráficos facilmente acessando as matrizes e vetores do PyCASSO, como pode ser visto no programa na Figura 2.4 e sua imagem gerada (Figura 2.5). Os cálculos necessários para o PCA e também para a Tomografia PCA se tornam simples contas usando pacotes matemáticos de python.

Gerar perfis radiais ou axiais também são de fundamental importância para o estudo de diversas propriedades galáticas. Podemos ver na Figura 2.6 (e em sua imagem gerada 2.7) um exemplo de perfil radia executado através da função `radialProfile`.

Dentro de nossa colaboração já existem mais de dez pessoas utilizando a biblioteca Py-

<sup>4</sup><http://matplotlib.org>

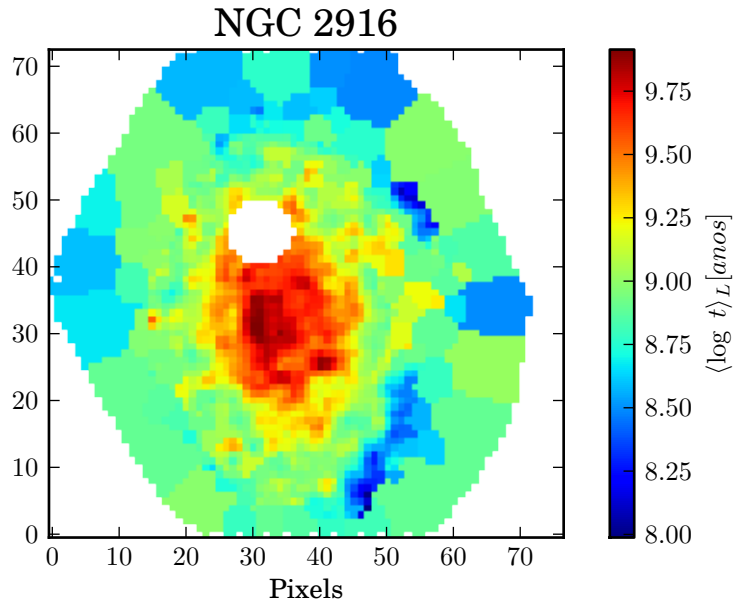
```
1 # Carregar arquivo FITS com os dados.
2 from pycasso import fitsQ3DataCube
3 K = fitsQ3DataCube('K0277_synthesis_suffix.fits')
4
5 # Acessar a idade media ponderada pela luminosidade.
6 at = K.at_flux__z
7
8 # Calcular a idade media da galaxia.
9 at_total = (at * K.Lobn__z).sum() / K.Lobn__z.sum()
10 print 'Idade media da galaxia K0277: \%.2f' \% at_total
```

**Figura 2.3:** Exemplo de acesso aos cubo de dados por arquivo FITS e o cálculo da idade estelar média de uma galáxia.

```
1 # Carregar arquivo FITS com os dados.
2 from pycasso import fitsQ3DataCube
3 K = fitsQ3DataCube('K0277_synthesis_suffix.fits')
4
5 # Converter zonas para imagem.
6 at_image = K.zoneToYX(K.at_flux__z, extensive=False)
7
8 # Desenhar o mapa.
9 import matplotlib.pyplot as plt
10 plt.imshow(at_image, origin='lower', interpolation='nearest')
11 plt.xlabel('Pixels')
12 cb = plt.colorbar()
13 cb.set_label(r'\langle \log t \rangle_L [anos]')
14 plt.title(r'\langle \log t \rangle_{L z}')
```

**Figura 2.4:** Programa para desenhar o mapa de idade estelar média ponderada pela luminosidade.





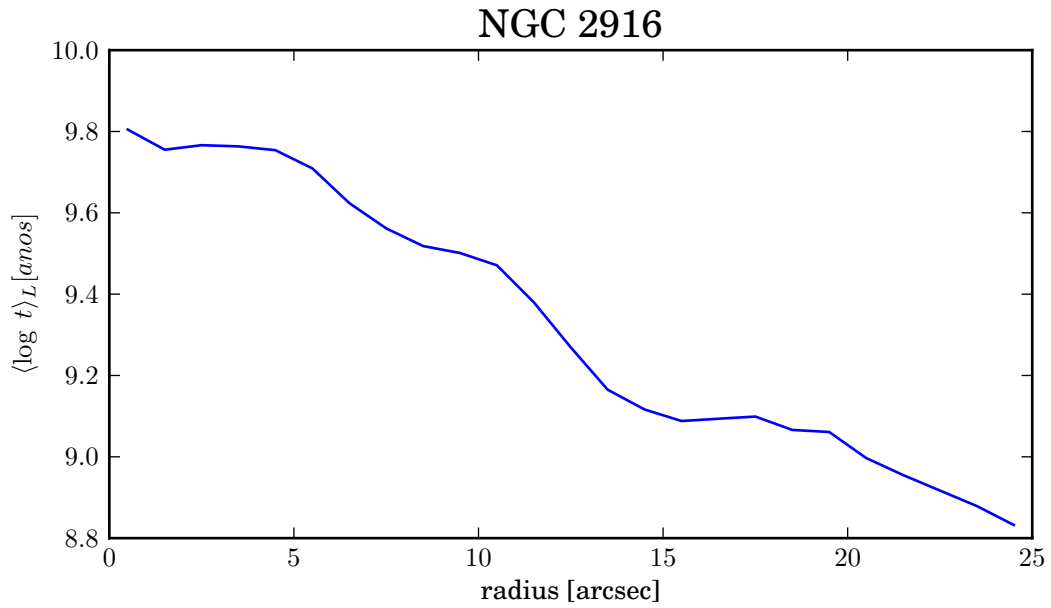
**Figura 2.5:** Mapa de idade estelar média ponderada pela luminosidade da galáxia NGC 2916 (CALIFA 277) gerado pelo programa da Figura 2.4.

```

1  # Carregar arquivo FITS com os dados.
2  from pycasso import fitsQ3DataCube
3  K = fitsQ3DataCube('K0277_synthesis_suffix.fits')
4
5  # Converter zonas para imagem.
6  at_image = K.zoneToYX(K.at_flux__z, extensive=False)
7
8  # Calcular o perfil radial.
9  bins = np.arange(0, 26, 1)
10 bin_center = (bins[1:] + bins[:-1]) / 2.0
11 at_rad = K.radialProfile(at_image, bins, rad_scale = 1.0)
12
13 # Desenhar o perfil.
14 import matplotlib.pyplot as plt
15 plt.xlabel('radius [arcsec]')
16 plt.ylabel(r'\langle \log t \rangle_L [anos]')
17 cb = plt.colorbar()
18 plt.plot(bin_center, at_rad)

```

**Figura 2.6:** Programa para desenhar o perfil radial da idade estelar média ponderada pela luminosidade.




**Figura 2.7:** Perfil radial da idade estelar média ponderada pela luminosidade da galáxia NGC 2916 (CALIFA 277) gerado pelo programa da Figura 2.6.

CASSO, com alguns artigos já publicados (Cid Fernandes et al. 2013, 2014; Pérez et al. 2013; González Delgado et al. 2013). Também há alguns usando indiretamente (Husemann et al. 2013; Iglesias-Páramo et al. 2013), uma tese de doutorado em curso e o presente trabalho.

## Capítulo 3

# PCA e Tomografia PCA

De medidas fisiológicas, como pulsação e respiração, até reconhecimento de padrões em sistemas complexos como reconhecimento facial e criptografia, passando por compactação de imagem, neurosciência e redução de ruídos em dados, podemos ver atuação de técnicas de PCA.

 Brincando com a própria técnica, o PCA em si é uma PC de um universo matemático-estatístico.

### 3.1 Principal Component Analysis

Baseada em encontrar os eixos com maiores variâncias em um conjunto de variáveis (no nosso caso, fluxos por lambda e por zona), a técnica PCA vem sendo de grande utilidade quando o assunto é estatística com muitas variáveis. Através de operações relativamente simples computacionalmente usando álgebra linear, é feito uma rotação na base de dados original, gerando uma nova base ortonormal não correlacionada através de um conjunto de autovalores e autovetores.

Existem diversas formas de se calcular essa base final. A prova matemática que você pode obter essa base é feita através de multiplicadores de Lagrange, calculando os autovetores e autovalores ( $\mathbf{e}_k$  e  $\lambda_k$ ) que maximizam o valor de  $\mathbf{e}_k^T \cdot \mathbf{C}_{cov} \cdot \mathbf{e}_k$  (ver Eq. 3.2) sujeito a restrição de que um autovetor deve ser ortogonal a qualquer outro da base ( $\mathbf{e}_i^T \mathbf{e}_j = 0$ ) e que todos devem ser normalizados ( $\mathbf{e}_i^T \mathbf{e}_i = 1$ ) (Jolliffe 2002, p. 5-6). No caso de PCA com espectros, cada um pode ser representado como um ponto num espaço  $\lambda$ -dimensional. Vários espectros formam uma nuvem de pontos. Assim, encontramos quais são os eixos mais significativos desse espaço em relação à variância, sujeitos às restrições acima. Fazemos esse cálculo encontrando

os autovetores e autovalores da matriz de correlação desse espaço (a matriz tem dimensão  $\lambda$ ) usando a biblioteca científica SciPy<sup>1</sup> (3.1), que encontra os autovetores e autovalores da matriz de correlação.

### 3.1.1 PCA das galáxias do CALIFA

Conforme a Seção 2.2 vimos que o cubo de espectros das galáxias do CALIFA estão acessíveis no PyCASSO separados por zonas. Os espectros observados estão armazenados em forma de uma matriz ( $n \times m$ ) com  $n$  zonas e  $m$  comprimentos de onda (`f_obs`<sup>2</sup> no PyCASSO).

$$\mathbf{F}_{z\lambda} = \begin{bmatrix} f_{z_0\lambda_0} & f_{z_0\lambda_1} & f_{z_0\lambda_2} & \cdots & f_{z_0\lambda_m} \\ f_{z_1\lambda_0} & f_{z_1\lambda_1} & f_{z_1\lambda_2} & \cdots & f_{z_1\lambda_m} \\ f_{z_2\lambda_0} & f_{z_2\lambda_1} & f_{z_2\lambda_2} & \cdots & f_{z_2\lambda_m} \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ f_{z_n\lambda_0} & f_{z_n\lambda_1} & f_{z_n\lambda_2} & \cdots & f_{z_n\lambda_m} \end{bmatrix} \quad (3.1)$$

Calculamos então o espectro médio de uma galáxia através da equação  $\langle \mathbf{F}_\lambda \rangle = (1/n) \sum_{i=0}^n f_{z_i\lambda}$  e então subtraímos a média de todos os espectros ( $\mathbf{I}_{z\lambda} = \mathbf{F}_{z\lambda} - \langle \mathbf{F}_\lambda \rangle$ ) para o cálculo da matriz de covariâncias usando um conjunto de dados com média zero. Vemos que a matriz de covariância possui dimensão  $\lambda \times \lambda$ .

$$\mathbf{C}_{cov} = \frac{[\mathbf{I}_{z\lambda}]^T \cdot \mathbf{I}_{z\lambda}}{n - 1} \quad (3.2)$$

Agora calculamos os autovalores e autovetores da matriz de covariância. Neste trabalho usamos o nome autoespectro para designar esses autovetores pois são de uma matriz de covariâncias entre espectros de cada zona. Então ordenamos-os decrescentemente pelo valor de seus autovalores. Os autoespectros são as PCs e os autovalores as variâncias. Isso feito, temos então o que necessitamos para iniciar o cálculo do Tomograma PCA. Muitas figuras de PCs e suas utilizações e diferenças nos pré-processamentos serão mostradas nos Capítulos 4 e 5, juntamente com a Tomografia PCA e as comparações com os parâmetros físicos da síntese de populações estelar.

<sup>1</sup><http://scipy.org/>

<sup>2</sup>No PyCASSO está amostrado de forma transposta ( $\lambda \times z$ ) a usada.

```
1 # Carregar arquivo FITS com os dados.
2 from pycasso import fitsQ3DataCube
3 K = fitsQ3DataCube('K0277_synthesis_suffix.fits')
4
5 # Calcular o espectro medio de uma galaxia.
6 # K.f_obs tem dimensao (lambda, zona), portanto,
7 # fazemos o espectro medio de todas as zonas.
8 f_obs_mean__l = K.f_obs.mean(axis = 1)
9
10 # Subtraimos a media
11 I_obs__zl = K.f_obs.transpose() - f_obs_mean__l
12
13 # Calcular a matrix de convariancia
14 import scipy as sp
15 n = K.N_zone
16 dot_product = sp.dot(I_obs__zl.transpose(), I_obs__zl)
17 covMat__ll = dot_product / (n - 1.0)
18
19 # Calcular os autovalores e autovetores
20 w, e = linalg.eigh(covMat__ll)
21
22 # Ordenar os autovetores decrescentemente pelo seu autovalor
23 S = sp.argsort(w)[::-1]
24 eigval = W[S]
25 eigvect = e[:, S]
```

**Figura 3.1:** Cálculo do procedimento completo de PCA para os espectros observados de uma galáxia do CALIFA usando o PyCASSO e a biblioteca científica de Python chamada SciPy. No final do código temos eigval e eigvect que são os autovalores e autovetores ordenados em forma decrescente.

## 3.2 Tomografia PCA

Na técnica de PCA procuramos os autoespectros (PCs) da matriz de correlação, ordenados pela variância, formam uma base ( $\mathbf{E}_{\lambda k}$ ) onde podemos projetar os nossos dados através da transformação:

$$\mathbf{T}_{zk} = \mathbf{I}_{z\lambda} \cdot \mathbf{E}_{\lambda k} \quad (3.3)$$

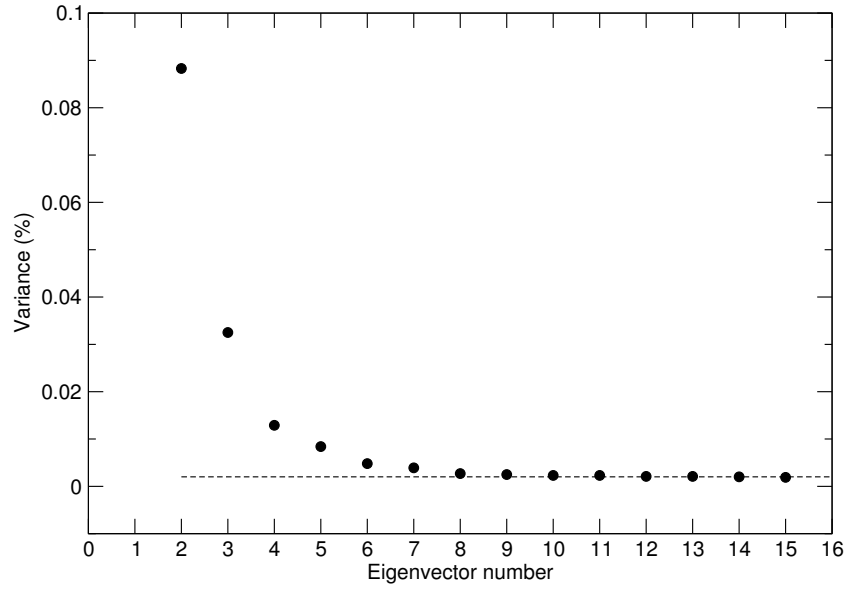
Projetamos nossa matriz de observáveis com a média subtraída ( $\mathbf{I}_{z\lambda}$ ) na base das PCs ( $\mathbf{E}_{\lambda k}$ ). Na posse dessa nova matriz transformada e de um mapa que leve de zona para uma par de coordenada ( $z \rightarrow (x, y)$ ), podemos montar assim uma imagem. Cada imagem funciona como uma “fatia” de um cubo de dados expandido na nova base, assim formando a Tomografia PCA<sup>3</sup>, criada e assim batizada por Steiner et al. (2009), que em seu artigo faz um paralelo com fatias de um espaço tridimensional (tomograma do corpo humano, por exemplo) ou no espaço de velocidades (Tomografia Doppler). Cada “fatia” possui um autoespectro relacionado que, em conjunto, trazem novas perspectivas e ideias para a interpretação de ambos. A passagem de coordenadas  $z \rightarrow (x, y)$  está exemplificada na Figura 2.5 através da função de zoneToYX dentro do PyCASSO.

### 3.2.1 Evidências de linhas largas

No artigo citado anteriormente, através do estudo dos autoespectros, e suas respectivas imagens, da galáxia LINER (*Low Ionizations Nuclear Emission-line Region*) NGC 4736, foram encontradas evidências de *broad lines* nas regiões nucleares da galáxia. Quando temos uma fonte que é capaz de produzir linhas largas no espectro é sinal da existência de um SMBH (*Super Massive Black Hole*). Cid Fernandes et al. (2004) mostra que a subtração bem detalhada das populações estelares nos espectros ajuda a encontrar linhas largas mais fracas em Seyferts-II, que são aquelas que possuem linhas estreitas, ajudando assim em na classificação desses objetos como tipo I ou II. O PCA, juntamente com a Tomografia PCA, fazem esse papel da subtração das populações estelares sem haver nenhuma parametrização.

Foram estudados os primeiros 8 autoespectros escolhidos através de um *scree test* (Figura 3.2), no qual se verifica a variância de cada PC e toma-se as mais relevantes. O autoespectro com mais variância (no artigo tratado com E1) possui 99.74% da variância e reproduz comportamento do gás e da população estelar somados (Figura 3.3). O segundo contribui com

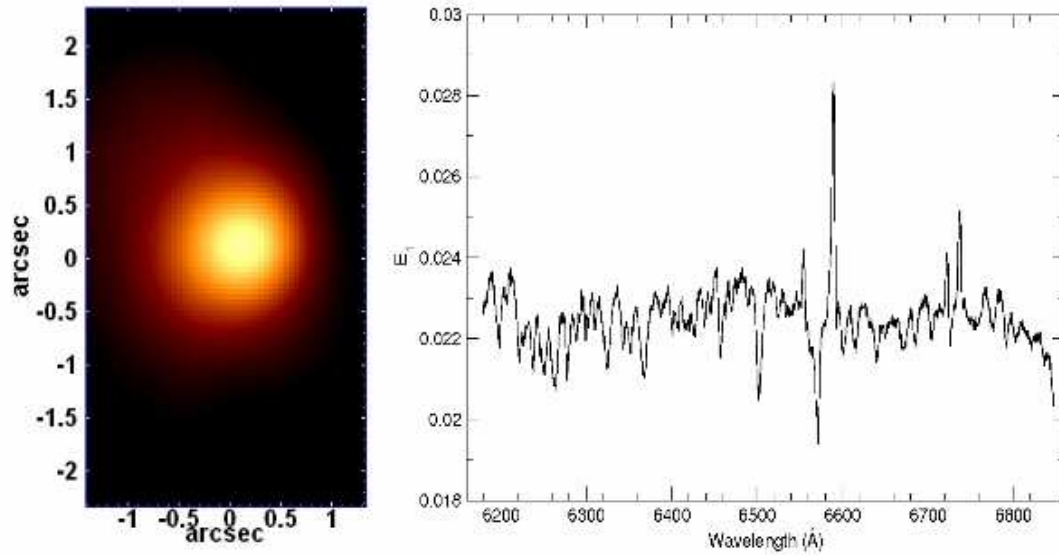
<sup>3</sup><http://www.astro.iag.usp.br/~pcatomography/>



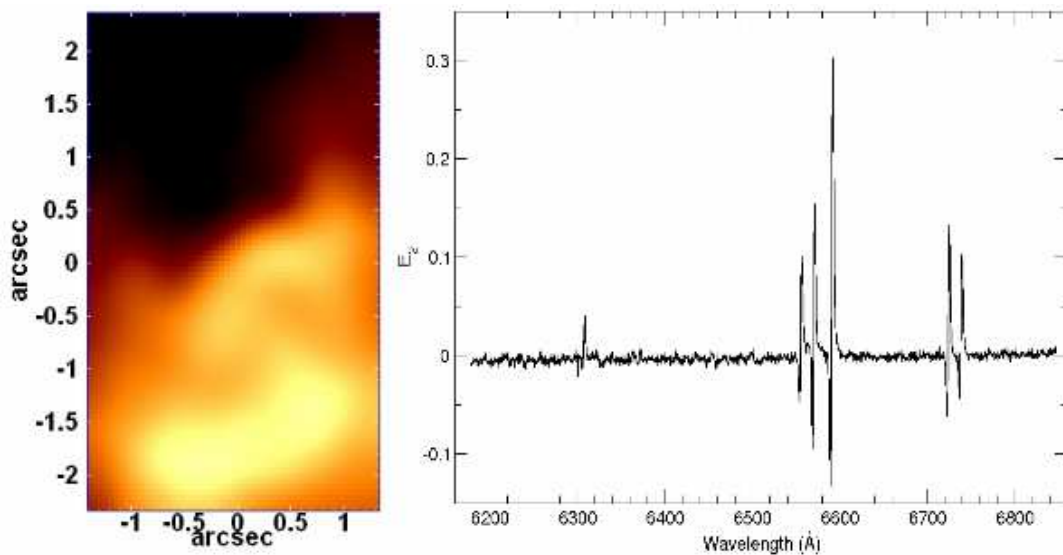
**Figura 3.2:** Scree test das primeiras 16 PCs do cubo de espectros da região central da galáxia NGC 4736. Retirado de Steiner et al. (2009, fig. 1).

0.088% para a variância e tem um claro padrão de rotação, tanto nas linhas do autoespectro quanto na imagem da tomografia (Figura 3.4). Mas é no terceiro (0.032% da variância) que mostra evidências de uma emissão larga de  $H\alpha$  (Figura 3.5). Essa assinatura é uma evidência típica de objetos com um AGN associado a um SMBH (Seyfert-I).

Com isso em mãos podemos estudar a Tomografia PCA nas galáxias do CALIFA. Temos disponíveis várias galáxias com os cubos de dados COMBO e aplicaremos essa técnica às galáxias presentes na Tabela 2.1 no Capítulo 5.

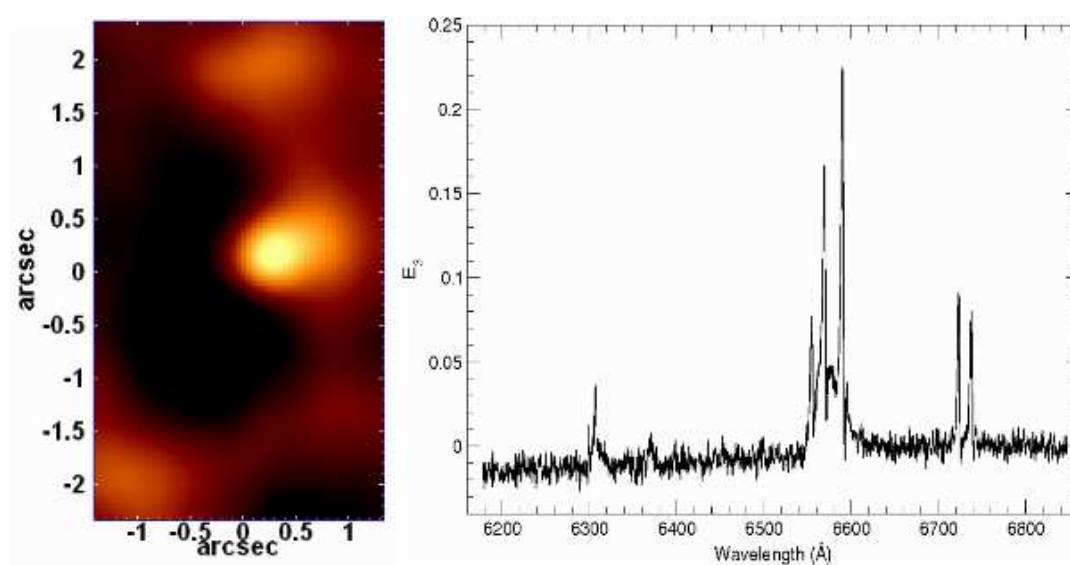


**Figura 3.3:** Autoespectro 1 e seu respectivo tomograma. Retirado de Steiner et al. (2009, fig. A1).



**Figura 3.4:** Autoespectro 2 e seu respectivo tomograma. Retirado de Steiner et al. (2009, fig. A2).





**Figura 3.5:** Autoespectro 2 e seu respectivo tomograma. Retirado de Steiner et al. (2009, fig. A3).

# Capítulo 4

## PCA nos espectros

O uso de métodos estatísticos já se estende por séculos em praticamente todas (senão todas) as áreas de conhecimento. Esse fato cria uma necessidade de que existam cada vez mais estudos sobre estudos, ou *metaestudos*<sup>1</sup>. Precisamos saber de que maneira os pré-processamentos de nossa amostra afetam os dados e, principalmente, o resultado após a aplicação de determinada técnica, para que dessa forma o desenvolvimento não se torne uma “caixa preta” inacessível.

### 4.1 Pré-processamento dos cubos

Antes dos espectros chegarem ao PyCASSO, todas as informações de *flags*<sup>2</sup> em *bad pixels* e linhas telúricas<sup>3</sup> são criadas em um *pipeline* de pré-processamento chamado QWICK. Esse *pipeline* também prepara os cubos para a execução do STARLIGHT, e para organização deles pelo PyCASSO, definindo as zonas de Voronoi, a reamostragem em  $\lambda$  e colocando os espectros em repouso usando o *redshift* calculado dentro dos 5” centrais da galáxia. Todas essas informações e pré-processamentos provenientes do QWICK são herdadas pelo PyCASSO e já estão contidas em seus cubos de espectros.

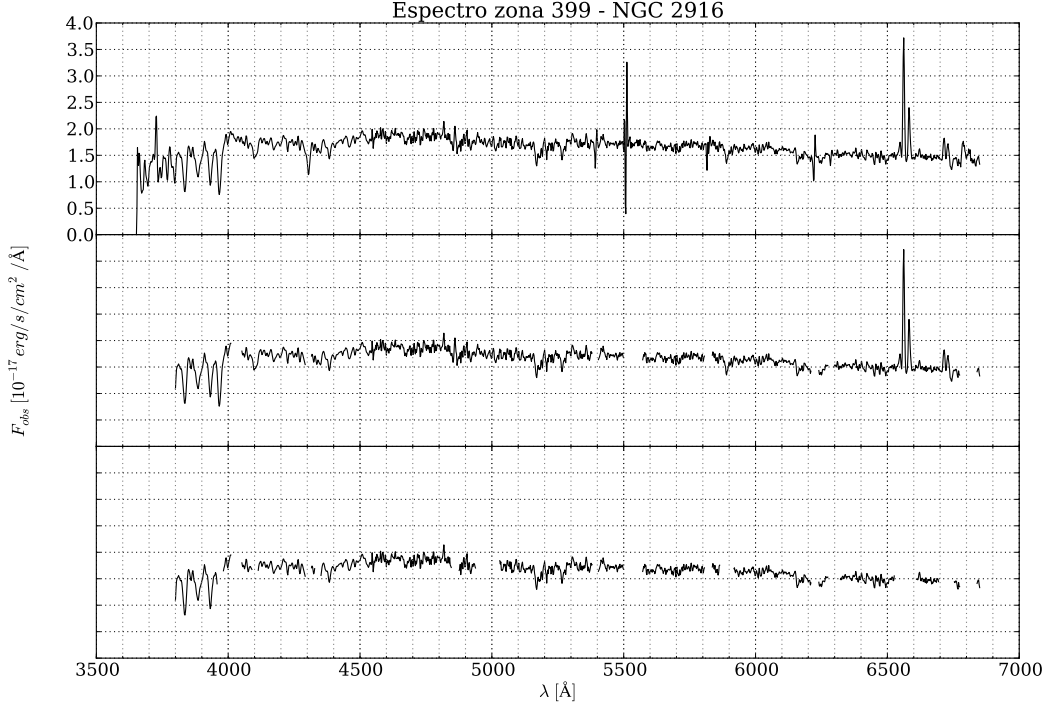
Apesar desses pré-processamentos supracitados, não é sobre eles vamos falar aqui, e sim sobre aqueles que são feitos nos espectros contidos no PyCASSO antes da aplicação do PCA. Como o PCA é uma técnica em qual calcula-se os eixos que, através da variância, melhor expandem sua base de dados, é natural que qualquer pré-processamento que altere a variância dos dados, resultará num conjunto diferente de PCs. Quando aplicamos o PCA aos cubos do

---

<sup>1</sup>Em alusão a metadados, que são dados sobre dados.

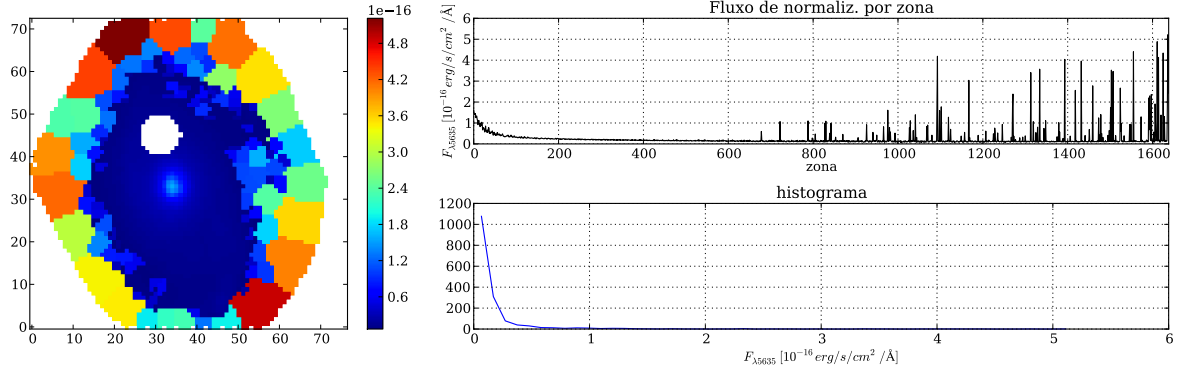
<sup>2</sup>Marcações.

<sup>3</sup>Linhas de absorção referentes à atmosfera.



**Figura 4.1:** Espectro da zona 399 da galáxia NGC 2916 (CALIFA 277). Acima está o espectro completo. No segundo vemos o espectro com linhas telúricas e bad-pixels removidos, além do limite de intervalo em comprimento de onda de 3800 a 6850 Å. No espectro mais abaixo, além das partes removidas no segundo, estão também removidas as mesmas linhas de emissão mascaradas na síntese de populações estelares.

CALIFA estamos buscando variâncias espaciais nos espectros. Durante nossas investigações fizemos uma série de testes com pré-processamentos nos espectros. Dois deles são constantes em todos os estudos. Primeiro todos os espectros são limitados ao intervalo de 3800 a 6850 Å. Após essa limitação, fazemos uma estatística com todos os *bad pixels* e linhas telúricas de cada cubo e aqueles que estão presentes em 95% dos espectros são removidos. Cabe aqui lembrar que todos os espectros precisam ter os mesmos pontos em  $\lambda$  pois precisamos construir a matriz de covariância. Podemos ver o efeito desses pré-processamentos nos espectros através dos 2 espectros mais acima, na Figura 4.1.



**Figura 4.2:** A imagem à esquerda e a superior direita, mostram o fluxo usado para a normalização de cada espectro em cada pixel (à esquerda) ou zona (superior direita). Na imagem inferior direita temos um histograma para valores do fluxo de normalização.

#### 4.1.1 Normalização

Imagine uma galáxia composta inteiramente pela mesma população estelar, em repouso, distribuídas da mesma maneira no espaço. Ou seja, em qualquer ponto da galáxia o espectro é o mesmo. Um PCA nessa galáxia hipotética nos mostraria apenas uma componente relevante. Adicione então a essa galáxia uma função de densidade de massa em função do raio, permitindo que de uma posição para outra a quantidade dessa determinada população se altera (mudando o brilho superficial da região). Uma componente nova irá surgir na sua análise PCA, mostrando que existe uma variância agora numa componente de escala (amplitude) nos espectros. Mas o que essa componente de escala nos diz sobre a física da população estelar existente? Essa componente seria um desperdício de variância para uma análise de populações estelares de uma galáxia.

No caso do CALIFA, o cubo de espectros inclui um  $FoV$  que abrange praticamente toda a galáxia. Isso gera uma variância indevida entres as zonas devido a luminosidade mais intensa nas zonas centrais da galáxia em comparação com as mais afastadas. Indevidas pois não trazem informação nova para a nossa análise, essas diferenças em amplitude não nos dizem nada sobre as populações estelares. Nas Figuras 4.3 e 4.4 vemos os quatro primeiros tomogramas para a galáxia NGC2916 sem e com normalização. Podemos notar que o primeiro autoespectro (e seu respectivo tomograma) no caso sem normalização mostra exatamente esse fator de escala. A PC1 se assemelha muito com a média (basta multiplicá-lo por -1). Seu tomograma, que reflete o peso dessa componente pra cada zona da galáxia, é claramente um fator de escala

(que pode ser considerado um fator de brilho). É fácil de notar a semelhança entre a primeira componente do caso com normalização e a segunda componente do caso sem normalização. Outra comparação que evidencia esse fator de escala é comparar a PC1 do caso sem normalização com os valores usados para a normalização dos espectros. Para a galáxia NGC2916 os valores estão na Figura 4.2. Nossas análises daqui para frente consideremos que cada espectro do cubo é normalizado pelo seu fluxo em 5635 Å.

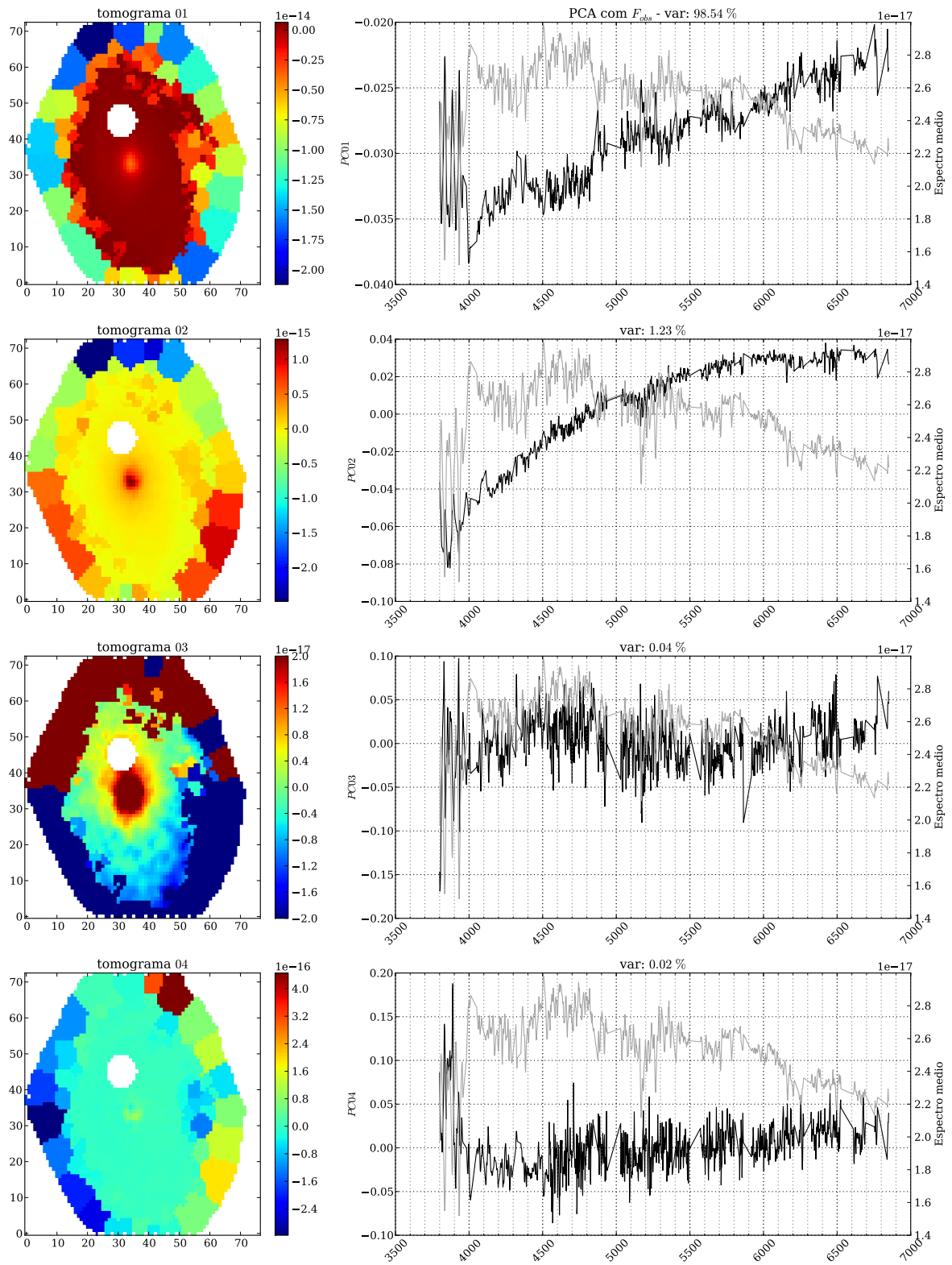
### 4.1.2 Cinemática

Usando novamente a ideia da galáxia hipotética com apenas uma população estelar, imagine agora que elas estão distribuídas uniformemente, mas estão em rotação com a galáxia. Da mesma forma o espectro de todas será igual salvo por deslocamentos em  $\lambda$  no espectro. Esses efeitos cinemáticos não estão nos trazendo informação alguma para o estudo das populações estelares e causam um grande desperdício de variâncias, sempre aparecendo nas primeiras PCs, além de existirem métodos mais eficazes e direcionados para a determinação de propriedades cinemáticas.

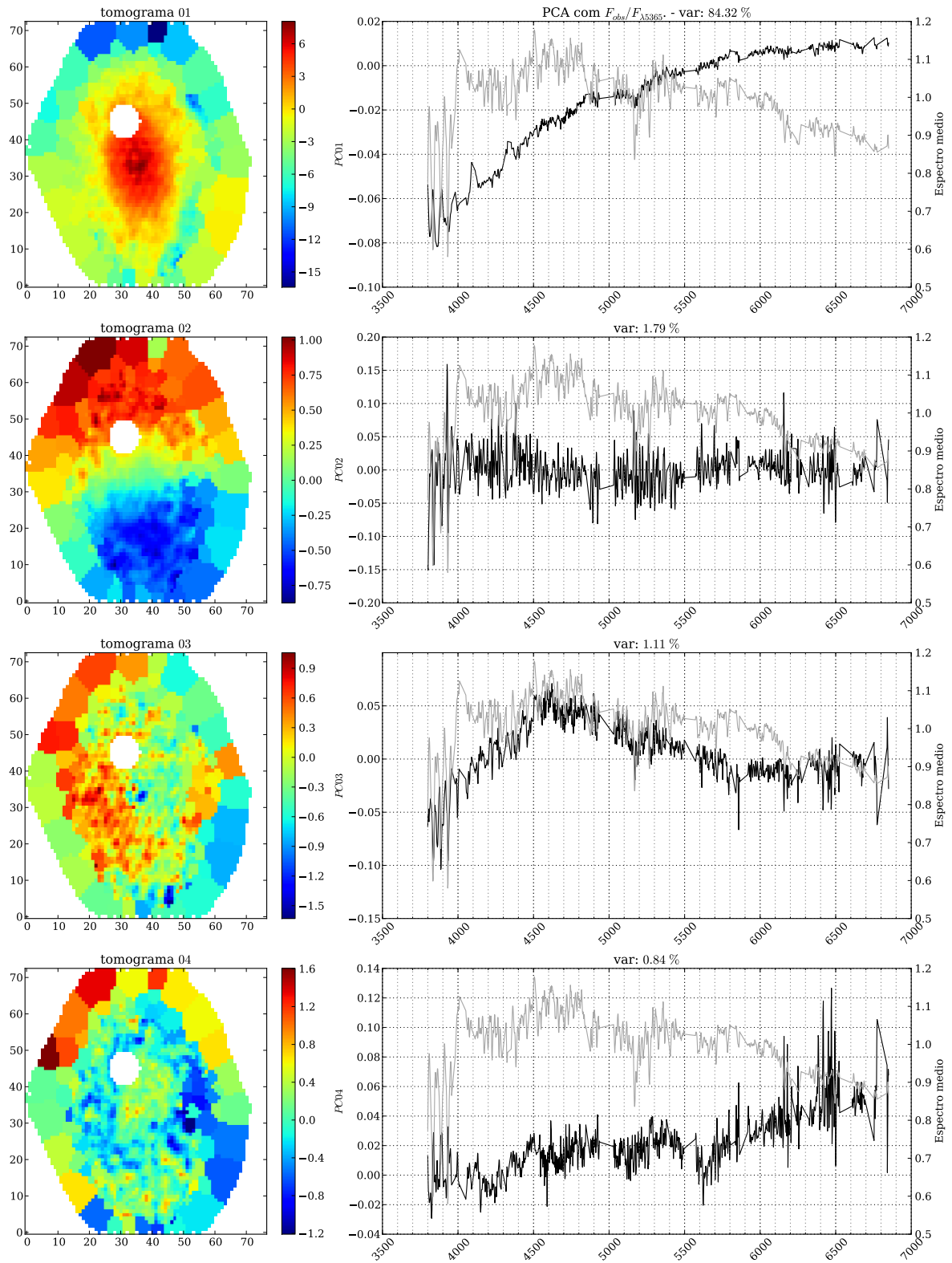
Pelo amplo *FoV* que as observações do CALIFA são feitas é normal que esse efeito apareça também, portanto os espectros aparecem com linhas deslocadas para o azul (*blue-shifted*) ou para o vermelho (*red-shifted*) dependendo da velocidade de rotação projetada. A dispersão de velocidades em cada ponto da galáxia também pode causar alargamento ou estreitamento das linhas. Podemos ver um padrão bem claro de rotação na PC2 da Figura 4.4. Da mesma forma, na PC3 da Figura 4.3. Nesta última, fizemos alguns ajustes na saturação das cores de modo que ficasse mais evidente o padrão de rotação, mas como não há a normalização ainda vemos algum efeito de intensidades misturado nessa PC, diferentemente da PC2 do PCA com normalização.

### 4.1.3 Linhas de emissão e intervalos específicos em comprimento de onda

Nos espectros, além dos *bad pixels* e linhas telúricas, podemos mascarar regiões desnecessárias para determinada investigação científica. Nosso foco é o estudo das populações estelares, portanto precisamos que as linhas de emissão, geralmente associadas ao gás presente nas galáxias [FIXME] sejam removidas do espectro. Dessa forma podemos fazer correlações entre os resultados do PCA e as propriedades físicas obtidas pela síntese.



**Figura 4.3:** Quatro primeiros PCs (e seus respectivos tomogramas) do PCA aplicado aos espectros sem normalização da galáxia NGC2916.



**Figura 4.4:** Quatro primeiros PCs (e seus respectivos tomogramas) do PCA aplicado aos espectros com normalização da galáxia NGC2916.

Podemos também executar o PCA apenas em intervalos específicos do espectro. Isso pode ser feito de várias formas, por exemplo utilizando todos os pontos do(s) intervalo(s) em  $\lambda$ , ou utilizando apenas o fluxo integrado, ou apenas larguras equivalentes de linhas, ou FWHM das linhas. Um exemplo aplicado é fazer o PCA apenas das regiões que abrangem  $[\text{O III}]/\text{H}\beta$  em conjunto com  $[\text{N II}]/\text{H}\alpha$  ou então  $\text{H}\delta$  e D4000 para estudar a variância espacial desses ratios.

Nas nossas análises no próximo capítulo 5 são mascaradas também todas as regiões que são removidas na síntese ( $\text{H}\epsilon$ : de 3960 a 3980 Å;  $\text{H}\delta$ : de 4092 a 4112 Å;  $\text{H}\gamma$ : de 4330 a 4350 Å;  $\text{H}\beta$ : de 4848 a 4874 Å;  $[\text{O III}]$ : de 4940 a 5028 Å;  $\text{He I}$  e  $\text{NaD}$ : de 5866 a 5916 Å;  $\text{H}\alpha$  e  $[\text{N II}]$ : de 6528 a 6608 Å;  $\text{S II}$ : de 6696 a 6752 Å).

#### 4.1.4 Fluxos observados e sintéticos

**[FIXME]** Isso não é um pré-processamento... aonde colocar???

Com o PyCASSO, temos o resultado da síntese de populações estelares já organizado para as galáxias do CALIFA. Com isso realizamos o PCA no cubo de espectros observados e no de espectros sintéticos, com e sem normalização. A grande diferença é que nos espectros sintéticos estão contidas apenas as informações sobre populações estelares<sup>4</sup>.

.....

**DAQUI PRA BAIXO AINDA É O ESQUELETO** !ojo! A Filosofia desse capítulo é aprender/testar como operar o PCA para que ele reflita isso ou aquilo... descrevemos uma série de experimentos nesse sentido.

Preprocessamentos e diferentes tipos de PCAs com ou sem linhas, diferentes faixas espectrais, com dados normalizados ou não. (importante!!!)

Vamos nos limitar a no-emission lines analysis, descrever a máscara de linhas de emissão, etc. Isso para facilitar a coisa e pq queremos correlar o resultado do PCA com os dados do Starlight (PyCASSO).

Simulações para ajudar a decifrar os resultados, população jovem + velha + modelo de distribuição espacial - ver efeitos de estratégias de preprocessamento.

Correlacionando os resultados do PCA com as propriedades do starlight (tipo de engenharia reversa)

<sup>4</sup>Suavização, correções por poeira e cinemática também são feitas nos espectros no processo de síntese. Mais detalhes em Cid Fernandes et al. (2005)



Linhas telúricas - remover ou não... bad pixels... mascarar ou não linhas de emissão

# Capítulo 5

## Resultados

 Aplicando PCA e Tomografia PCA para 3 galáxias espirais, 3 elípticas e 3 mergers.

## Capítulo 6

# Conclusões e perspectivas

### 6.1 Este trabalho

!ojo!

### 6.2 Trabalhos futuros

!ojo! Remover cinemática ( $v_0$  &  $v_d$ ) pois é um desperdício de variância. Incluir linhas de emissão na análise (aqui posso incluir umas figuras com essa parte que já está programada).

# Apêndice A

## **[FIXME]** Python



# Referências Bibliográficas

- Abazajian, K. N., Adelman-McCarthy, J. K., Agüeros, M. A., Allam, S. S., Allende Prieto, C., An, D., Anderson, K. S. J., Anderson, S. F. et al. 2009, *ApJS*, 182, 543
- Ahn, C. P., Alexandroff, R., Allende Prieto, C., Anderson, S. F., Anderton, T., Andrews, B. H., Aubourg, É., Bailey, S. et al. 2012, *ApJS*, 203, 21
- Balakrishnan, G., Durand, F., & Guttag, J. 2013, 2013 IEEE Conference on Computer Vision and Pattern Recognition, 0, 3430
- Benítez, N., Gaztañaga, E., Miquel, R., Castander, F., Moles, M., Crocce, M., Fernández-Soto, A., Fosalba, P. et al. 2009, *ApJ*, 691, 241
- Blanc, G. A., Gebhardt, K., Heiderman, A., Evans, II, N. J., Jogee, S., van den Bosch, R., Marinova, I., Weinzierl, T. et al. 2010, in *Astronomical Society of the Pacific Conference Series*, Vol. 432, *New Horizons in Astronomy: Frank N. Bash Symposium 2009*, ed. L. M. Stanford, J. D. Green, L. Hao, & Y. Mao, 180
- Borcea, L., Callaghan, T., & Papanicolaou, G. 2012, *CoRR*, abs/1208.3700
- Chen, Y.-M., Kauffmann, G., Tremonti, C. A., White, S., Heckman, T. M., Kovač, K., Bundy, K., Chisholm, J. et al. 2012, *MNRAS*, 421, 314
- Cid Fernandes, R., González Delgado, R. M., García Benito, R., Pérez, E., de Amorim, A. L., Sánchez, S. F., Husemann, B., Falcón Barroso, J. et al. 2014, *A&A*, 561, A130
- Cid Fernandes, R., González Delgado, R. M., Schmitt, H., Storchi-Bergmann, T., Martins, L. P., Pérez, E., Heckman, T., Leitherer, C. et al. 2004, *ApJ*, 605, 105
- Cid Fernandes, R., Mateus, A., Sodré, L., Stasińska, G., & Gomes, J. M. 2005, *MNRAS*, 358, 363
- Cid Fernandes, R., Pérez, E., García Benito, R., González Delgado, R. M., de Amorim, A. L., Sánchez, S. F., Husemann, B., Falcón Barroso, J. et al. 2013, *A&A*, 557, A86
- Colless, M. 1999, in *Large-Scale Structure in the Universe*, ed. G. Efstathiou & et al., 105
- Croom, S. M., Lawrence, J. S., Bland-Hawthorn, J., Bryant, J. J., Fogarty, L., Richards, S., Goodwin, M., Farrell, T. et al. 2012, *MNRAS*, 421, 872

- da Costa, L. N., Pellegrini, P. S., Sargent, W. L. W., Tonry, J., Davis, M., Meiksin, A., Latham, D. W., Menzies, J. W. et al. 1988, *ApJ*, 327, 544
- de Zeeuw, P. T., Bureau, M., Emsellem, E., Bacon, R., Carollo, C. M., Copin, Y., Davies, R. L., Kuntschner, H. et al. 2002, *MNRAS*, 329, 513
- González Delgado, R. M., Pérez, E., Cid Fernandes, R., García-Benito, R., de Amorim, A. L., Sánchez, S. F., Husemann, B., Cortijo-Ferrero, C. et al. 2013, *ArXiv e-prints*
- Huchra, J., Davis, M., Latham, D., & Tonry, J. 1983, *ApJS*, 52, 89
- Huchra, J. P. 1988, in *Astronomical Society of the Pacific Conference Series*, Vol. 5, *The Minnesota lectures on Clusters of Galaxies and Large-Scale Structure*, ed. J. M. Dickey, 41–70
- Husemann, B., Jahnke, K., Sánchez, S. F., Barrado, D., Bekerait\*error\*è, S., Bomans, D. J., Castillo-Morales, A., Catalán-Torrecilla, C. et al. 2013, *A&A*, 549, A87
- Iglesias-Páramo, J., Vílchez, J. M., Galbany, L., Sánchez, S. F., Rosales-Ortega, F. F., Mast, D., García-Benito, R., Husemann, B. et al. 2013, *A&A*, 553, L7
- Ivezic, Z., Tyson, J. A., Acosta, E., Allsman, R., Anderson, S. F., Andrew, J., Angel, R., Axelrod, T. et al. 2008, *ArXiv e-prints*
- Jolliffe, I. 2002, *Principal Component Analysis*, 2nd edn., *Springer series in statistics* (Springer)
- Kamruzzaman, S. M., Siddiqi, F. A., Islam, M. S., Haque, M. E., & Alam, M. S. 2010, *CoRR*, abs/1009.4974
- Kelz, A., Verheijen, M. A. W., Roth, M. M., Bauer, S. M., Becker, T., Paschke, J., Popow, E., Sánchez, S. F. et al. 2006, *PASP*, 118, 129
- Mateus, A., Sodrè, L., Cid Fernandes, R., & Stasińska, G. 2007, *MNRAS*, 374, 1457
- Pérez, E., Cid Fernandes, R., González Delgado, R. M., García-Benito, R., Sánchez, S. F., Husemann, B., Mast, D., Rodón, J. R. et al. 2013, *ApJ*, 764, L1
- Ricci, T. V., Steiner, J. E., & Menezes, R. B. 2011, *ApJ*, 734, L10
- Riffel, R., Riffel, R. A., Ferrari, F., & Storchi-Bergmann, T. 2011, *MNRAS*, 416, 493
- Rosales-Ortega, F. F., Kennicutt, R. C., Sánchez, S. F., Díaz, A. I., Pasquali, A., Johnson, B. D., & Hao, C. N. 2010, *MNRAS*, 405, 735
- Roth, M. M., Kelz, A., Fechner, T., Hahn, T., Bauer, S.-M., Becker, T., Böhm, P., Christensen, L. et al. 2005, *PASP*, 117, 620
- Sánchez, S. F., Kennicutt, R. C., Gil de Paz, A., van de Ven, G., Vílchez, J. M., Wisotzki, L., Walcher, C. J., Mast, D. et al. 2012, *A&A*, 538, A8

- Skrutskie, M. F., Cutri, R. M., Stiening, R., Weinberg, M. D., Schneider, S., Carpenter, J. M., Beichman, C., Capps, R. et al. 2006, *AJ*, 131, 1163
- Steiner, J. E., Menezes, R. B., Ricci, T. V., & Oliveira, A. S. 2009, *MNRAS*, 395, 64
- Verheijen, M. A. W., Bershady, M. A., Andersen, D. R., Swaters, R. A., Westfall, K., Kelz, A., & Roth, M. M. 2004, *Astronomische Nachrichten*, 325, 151
- York, D. G., Adelman, J., Anderson, Jr., J. E., Anderson, S. F., Annis, J., Bahcall, N. A., Bakken, J. A., Barkhouser, R. et al. 2000, *AJ*, 120, 1579