Machine Learning techniques for Stellar Spectroscopy

Part 2

Part 1

- What is Machine Learning?
- How does it work?
- Brief introduction to Artificial Neuronal Networks.

E Millenter

Lorenzo Spina

INAF - Astronomical Observatory of Padua (Italy)

© Pete Ryan

spina.astro@gmail.com

@lorenzospina

Part 2

- The Cannon
- The Payne
- Issues of using ML for stellar spectroscopy and possible solutions

Moore's Law: The number of transistors on microchips doubles every two years Our World

Moore's law describes the empirical regularity that the number of transistors on integrated circuits doubles approximately every two years. This advancement is important for other aspects of technological progress in computing – such as processing speed or the price of computers.



OurWorldinData.org – Research and data to make progress against the world's largest problems.

Licensed under CC-BY by the authors Hannah Ritchie and Max Roser.

in Data



OurWorldinData.org – Research and data to make progress against the world's largest problems.

Licensed under CC-BY by the authors Hannah Ritchie and Max Roser.



OurWorldinData.org - Research and data to make progress against the world's largest problems.

Licensed under CC-BY by the authors Hannah Ritchie and Max Roser.



OurWorldinData.org – Research and data to make progress against the world's largest problems.

Licensed under CC-BY by the authors Hannah Ritchie and Max Roser.



In our dreams...



In our dreams...



In our dreams...



ML algorithm

(magics happen)







a cat



a cat



a translated cat: the same cat



a cat



a translated cat: the same cat



a distorted cat: the same cat



a cat



a translated cat: the same cat



a distorted cat: the same cat

a spectrum



a cat



a spectrum



a translated cat: the same cat



a distorted cat: the same cat



a translated spectrum: not the same spectrum!



a cat



a spectrum



a translated cat: the same cat



a distorted cat: the same cat



a translated spectrum: not the same spectrum!



a distorted spectrum: not the same spectrum!

The Cannon Ness et al. (2015)

A pure data-driven approach to Spectroscopic Analysis.

Flux at every pixel for the nth-star $f_{n\lambda} = g(\mathbf{l}_n | \theta_{\lambda}) +$ noise Stellar parameters and abundances of the nth-star. (known labels) Goefficients.

The **coefficients** and the **noise** can be "learnt" from a training set of spectra whose stellar parameters are already know. Once the model is ready, you can apply it to other stars with unknown labels.

The Cannon relies on the following to be true:

- stellar flux varies smoothly with stellar labels
- stars with the same labels have the same spectra

The Cannon Ness et al. (2015)

A pure data-driven approach to Spectroscopic Analysis.

Flux at every pixel for the nth-star $f_{n\lambda} = g(\mathbf{l}_n | \theta_{\lambda}) +$ noise Stellar parameters and abundances of the nth-star. (known labels) Goefficients.

The g function is a polynomium, whose order is set by the user.

First-order polynomium. The most simple model.

Authors suggest to use a second-order polynomium (quadratic-in-label case).

$$f_{n\lambda} = a_{\lambda} + b_{\lambda} (\text{Teff})_n + c_{\lambda} (\log g)_n + d_{\lambda} ([\text{Fe}/\text{H}])_n + \text{noise}$$
$$\theta_{\lambda} = (a_{\lambda}, b_{\lambda}, c_{\lambda}, d_{\lambda})$$

 $\begin{aligned} F_{n\lambda} &= \theta_0 \qquad (\text{constant term}) \\ &+ \theta_{T_{\text{eff}}} T_{\text{eff}} + \dots + \theta_{X_N} [X_N/\text{Fe}] \qquad (\text{linear terms}) \\ &+ \theta_{T_{\text{eff}}^2} T_{\text{eff}}^2 + \dots + \theta_{X_N^2} ([X_N/\text{Fe}])^2 \qquad (\text{squared terms}) \\ &+ \theta_{T_{\text{eff}} \log(g)} T_{\text{eff}} \log(g) + \dots \\ &+ \theta_{X_N X_{N-1}} [X_N/\text{Fe}] [X_{N-1}/\text{Fe}] \qquad (\text{cross-terms}) \end{aligned}$

 $+ \operatorname{error}$

The Cannon Ness et al. (2015)

A pure data-driven approach to Spectroscopic Analysis.



The Cannon Ness et al. (2015)

A pure data-driven approach to Spectroscopic Analysis.



The Cannon Ness et al. (2015)

A pure data-driven approach to Spectroscopic Analysis.



@lorenzospina





The Cannon Ness et al. (2015)

A pure data-driven approach to Spectroscopic Analysis.







Advantage

The method is extremely symple and fast.

The Cannon Ness et al. (2015)

A pure data-driven approach to Spectroscopic Analysis.



Advantage

The method is extremely symple and fast.

Danger

The method uses no physical stellar models.

Beware of hidden bias and correlations!

The danger of using datadriven models.

We need to be extra cautious of what the ML algorithms are "learning" and what training sets we build.

The effect of **biases** in the training set...

Take-home messages:

The training set must be representative of the specific problem we want to address. Any bias that exists in the training set will be folded into the test set.

The effect of **correlations** in the training set.

Is the model "learning" **measurements**?

...or is it "learning" hidden correlations?

The effect of **correlations** in the training set.

Beware of hidden bias and correlations!

The danger of using datadriven models.

We need to be extra cautious of what the ML algorithms are "learning" and what training sets we build.

Beware of hidden bias and correlations!

The danger of using datadriven models.

We need to be extra cautious of what the ML algorithms are "learning" and what training sets we build.

Unfortunately bias and correlations are unavoidable in stellar spectroscopy...

The Payne Ting et al. (2019)

Physical generalization: the training sample is composed by synthetic spectra.

The Payne Ting et al. (2019)

Physical generalization: the training sample is composed by synthetic spectra.

Full control on the training sample You can sample the entire space of parameters.

The Payne Ting et al. (2019)

Physical generalization: the training sample is composed by synthetic spectra.

Full control on the training sample You can sample the entire space of parameters.

Thanks to this fast interpolator of atmospheric models you can simply fit your spectra in order to calculate the atmospheric parameters.

The Payne Ting et al. (2019)

Physical generalization: the training sample is composed by synthetic spectra.

@lorenzospina

12

Models are good approximations of reality, but they are not perfect.

Synthetic gap: differences in feature distributions between synthetic and observed spectra.

Tips to reduce the synthetic gap:

- pre-processing the spectra: matching the resolution and sampling of the spectra to a common wavelength scale and removing the continuum.
- argumenting the spectra data set: adding Gaussian noise, rotational and radial velocities, masking telluric regions, and zeroing flux values to mimic bad pixels

Models are good approximations of reality, but they are not perfect.

Synthetic gap: differences in feature distributions between synthetic and observed spectra.

Tips to reduce the synthetic gap:

- pre-processing the spectra: matching the resolution and sampling of the spectra to a common wavelength scale and removing the continuum.
- argumenting the spectra data set: adding Gaussian noise, rotational and radial velocities, masking telluric regions, and zeroing flux values to mimic bad pixels

Before data argumentation After data argumentation

t-SNE visualization of the synthetic and APOGEE spectra

Models are good approximations of reality, but they are not perfect.

Synthetic gap: differences in feature distributions between synthetic and observed spectra.

Tips to reduce the synthetic gap:

- pre-processing the spectra: matching the resolution and sampling of the spectra to a common wavelength scale and removing the continuum.
- argumenting the spectra data set: adding Gaussian noise, rotational and radial velocities, masking telluric regions, and zeroing flux values to mimic bad pixels

t-SNE visualization of the synthetic and APOGEE spectra

Models are good approximations of reality, but they are not perfect. **Synthetic gap**: differences in feature distributions between synthetic ar

Tips to reduce the synthetic gap:

- pre-processing the spectra: matching the resolution and sampling of scale and removing the continuum.
- argumenting the spectra data set: adding Gaussian noise, rotational a regions, and zeroing flux values to mimic bad pixels

t-SNE visualization of the synthetic and APOGEE spectra

Fabbro et al. (2017)

Spectra generalisation

The two approaches with machine learning

selection functions and hiddden biases

physics.

The two approaches with machine learning

Cycle-Starnet O'Briain et al. (2021)

Cycle-Starnet O'Briain et al. (2021)

Cycle-Starnet

Limitations:

- It is not always straightforward to make the algorithm converge during training.
- It relies on the assumption that the uncertainties in the observed spectra are perfectly known. How a mis-characterisation of the noise affect the final results is yet to be studied.
- Labels from the synthetic and observad domains must span the same range. However, for a new survey we may not know the range of stellar labels.
- Cycle-Starnet may not extrapolate well for outliers and exotic stellar objects.
- Cycle-Starnet can correct for inaccurate atomic parameters in the linelist. However, this adds much more complexity in the analysis.
- It is impossible to correct for all the possible unknowns. The synthetic gap is reduced but not completely eliminated.

However, the study of Tranfer Features algorithms is currently proceeding at high speed. It is very likely that in the next few years we will have better algorithms that could alleviate all these limitations.

Conclusions

- Machine learning algorithms are a powerful tool.
 Significant computational power, huge amount of data, new algorithms: we can embrace new methods in our research, passing from a model-drive approach to a data-driven approach.
- Machine learning algorithms applied to stellar spectroscopy have some pitfalls and limitations.
 Data-driven models: bias and correlations.

Physical generalisation models: synthetic gap.

These are unavoidable pitfalls. However there are codes that achieve reasonably good results.

...also see Nandakumar et al. (2020) for a nice example of data homogeneisation (GALAH-APOGEE) with the Cannon.

• There is a huge effort in researching new machine learning algorithms. The "state-of-the-art" is fluid and can evolve very rapidly. Stay tuned!