# Introducción a la Clasificación Automática de Texto con PySS3
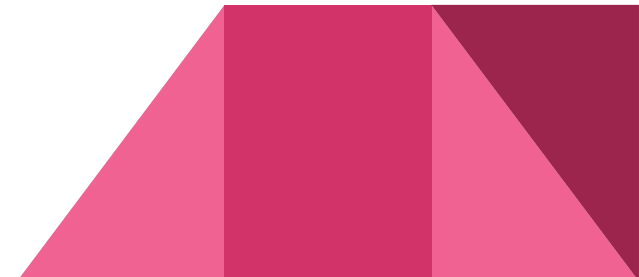
Sergio G. Burdisso[1,2]

1 Universidad Nacional de San Luis (UNSL), Argentina.
2 Consejo Nacional de Investigaciones Científicas y Técnicas (CONICET), Argentina.

# Outline

- **¿Qué es la Clasificación Automática de textos/documentos?**
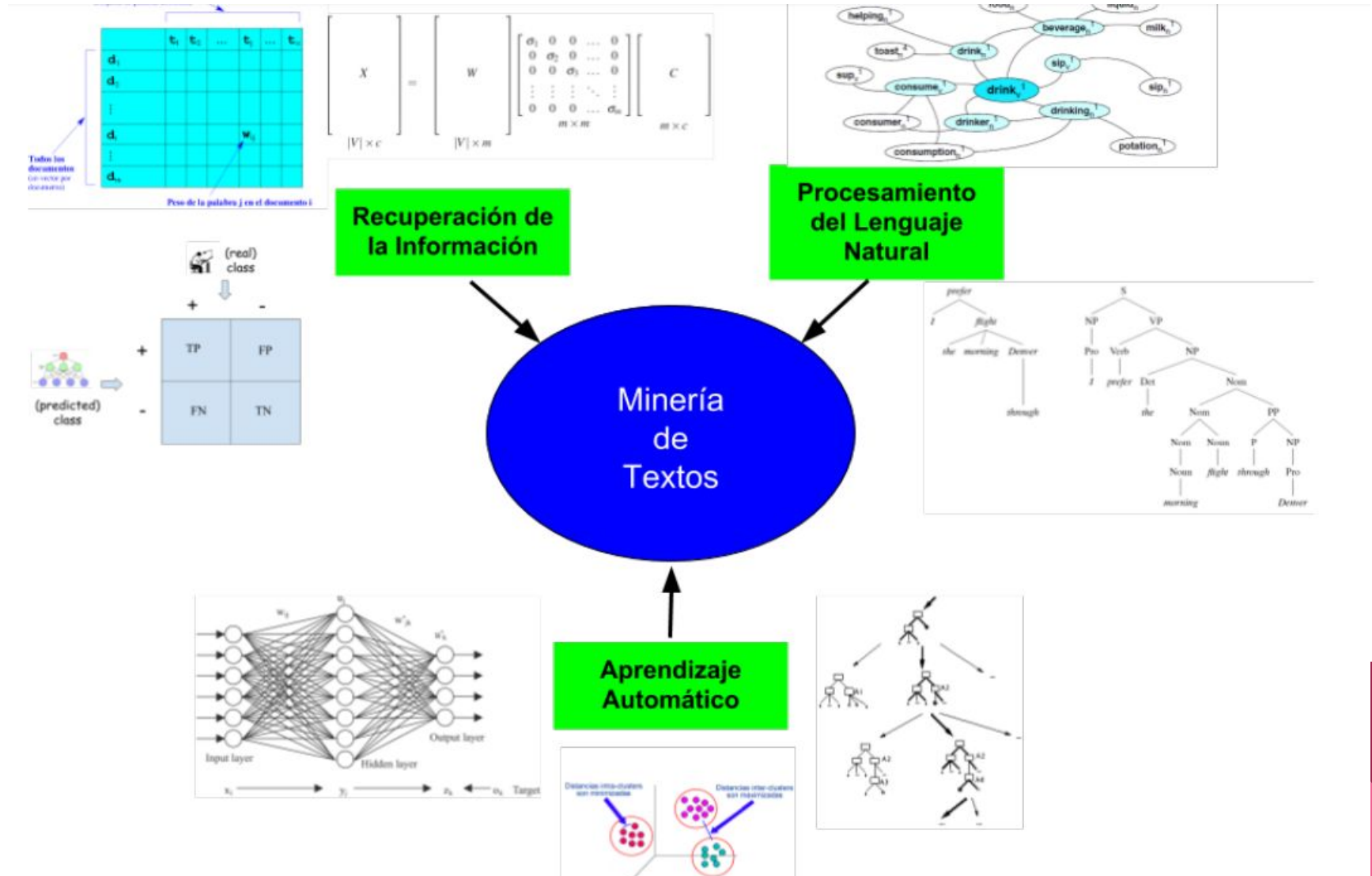
- **¿Qué es SS3?**

- **¿Qué es PySS3?**

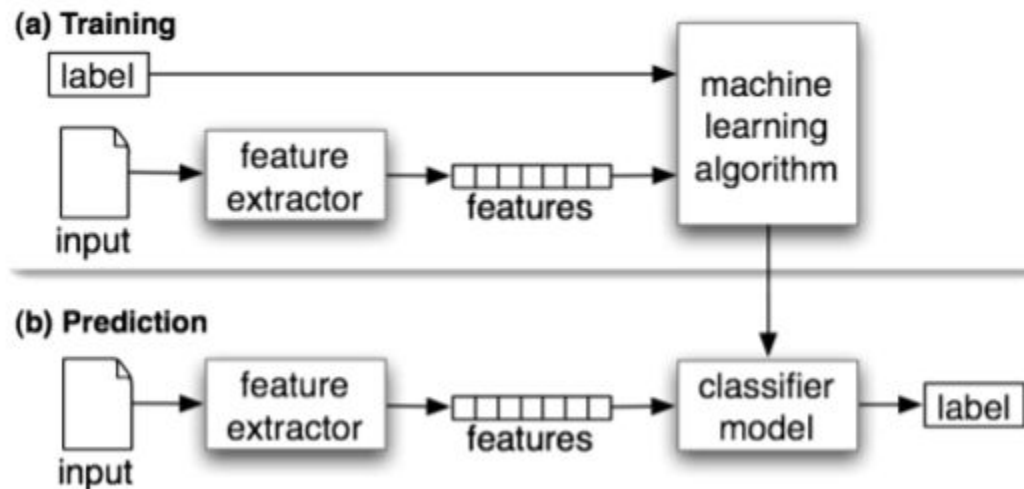# Clasificación Automática de Texto/Documentos

# Clasificación Automática de Texto

# Clasificación Automática de Texto (Contexto)

# Clasificación Automática de Texto (Detalles)

# Clasificación Automática de Texto (Ejemplos)

- Sentiment Analysis
  - **Social media monitoring:** analyze tweets and/or Facebook comments and detect if they are talking positively or negatively about a brand.
  - **Customer service:** analyze support queries to quickly detect angry and frustrated customers.
  - **Customer feedback:** analyze comments or survey responses to find if customers like or dislike particular aspects of a product or service.
- Language Detection
  - Routing customer support tickets to the correct team.
  - Sort through documents according to their languages.
  - Filtering incoming messages in undesired languages.
- Spam filtering
- Topic Categorization
- Profiling
- **Health and Safety**
  - **Early risk prediction on the Internet** (e.g. Depression, Anorexia, Self-harm, Terrorism, etc.)

# The SS3 classification model

Supervised Machine Learning Model for Text Classification
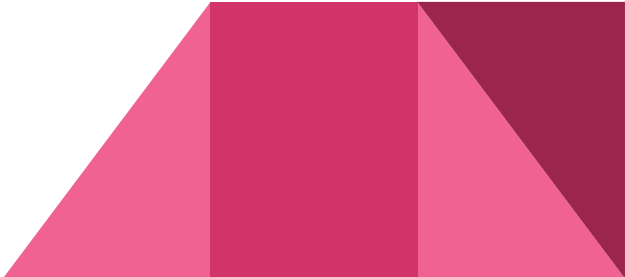
# The SS3 text classifier (1/3)

First, it uses a **function**, *gv*, to value **word relevance** relative to **each category**, for example:

$$gv(\text{'sushi'}, food) = 0.85; \qquad gv(\text{'the'}, food) = 0;$$
$$gv(\text{'sushi'}, music) = 0.09; \qquad gv(\text{'the'}, music) = 0;$$
$$gv(\text{'sushi'}, health) = 0.50; \qquad gv(\text{'the'}, health) = 0;$$
$$gv(\text{'sushi'}, sports) = 0.02; \qquad gv(\text{'the'}, sports) = 0;$$
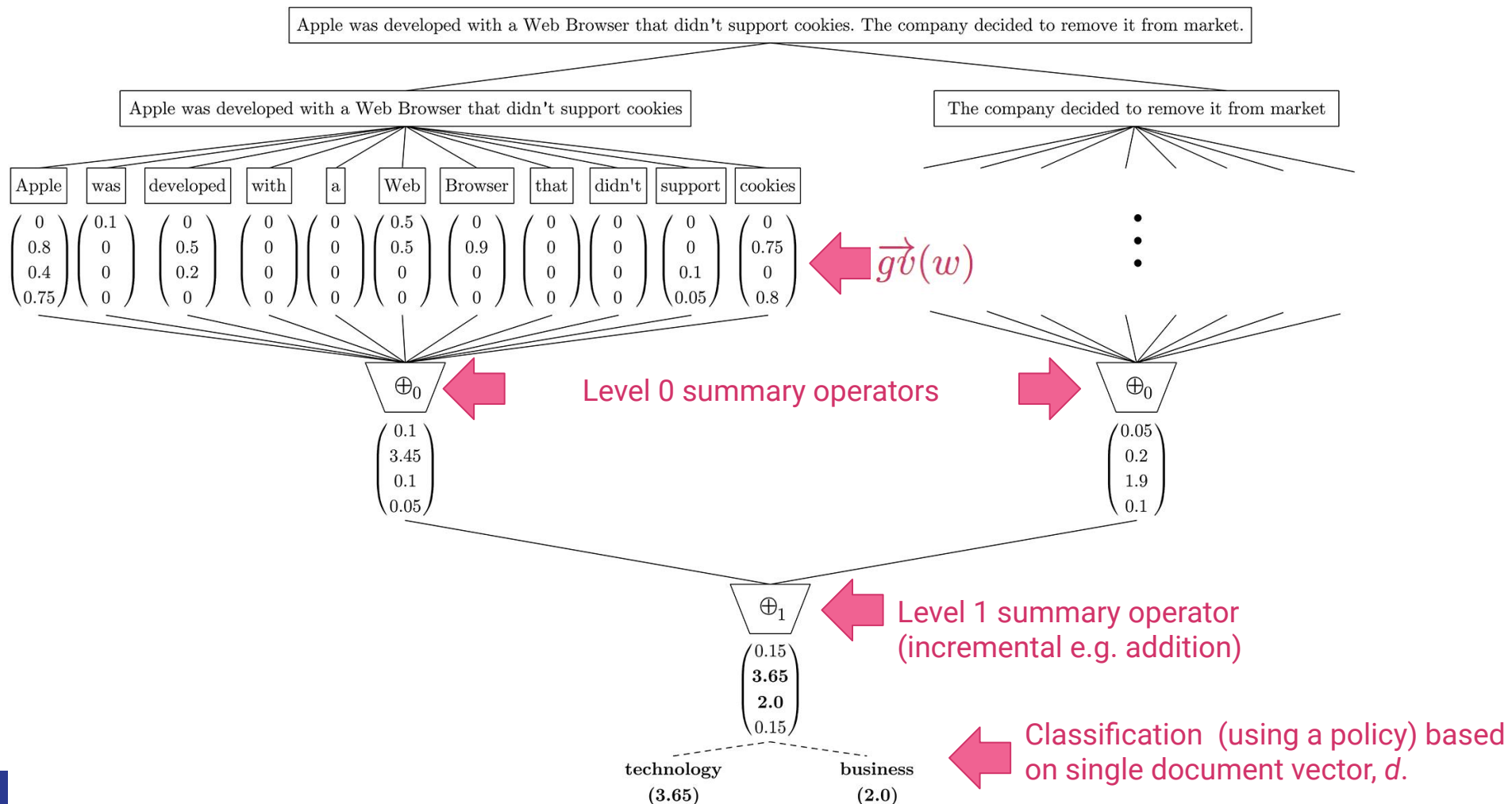
This function has a **vectorial version**:

$$\vec{gv}(w) = (gv(w, c_0), gv(w, c_1), \ldots, gv(w, c_k))$$

Thus, the previous example becomes:

$$\vec{gv}(\text{'sushi'}) = (0.85, 0.09, 0.5, 0.02);$$
$$\vec{gv}(\text{'the'}) = (0, 0, 0, 0);$$

# The SS3 text classifier (2/3)

Then, it converts the input into a **hierarchy of blocks**, computing a *confidence vector* for each hierarchy block. Classification is made using the single document *confidence vector*.
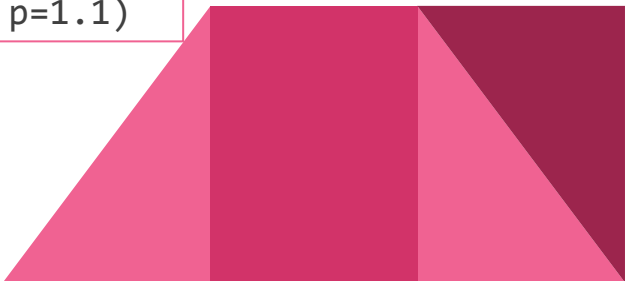


Level 0 summary operators

Level 1 summary operator
(incremental e.g. addition)

Classification (using a policy) based on single document vector, *d*.

# The SS3 text classifier (3/3) - Hyperparameters

$$global\ value = local\ value \cdot significance \cdot sanction$$

- σ ("Smoothness")
- λ ("Significance")
- ρ ("Sanction")

```
clf = SS3(s=0.32, l=1.24, p=1.1)
```
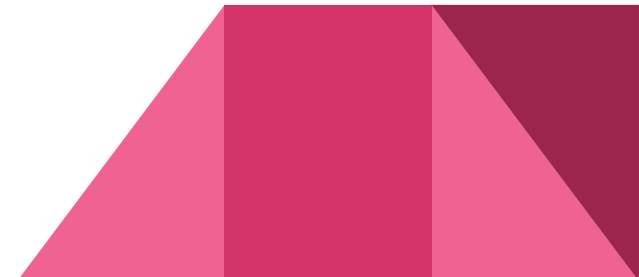
```
clf.set_hyperparameters(s=0.32, l=1.24, p=1.1)
```

# PySS3 - ¿Qué es?

Un paquete python que implementa SS3 y que además viene con un conjunto de herramientas de desarrollo y visualización.

Compuesto por:

- Módulo Principal
- 3 Submódulos:
  - pyss3.server
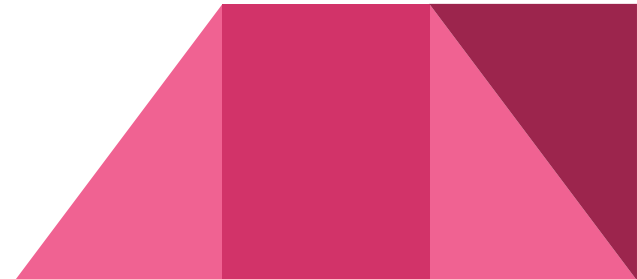  - pyss3.cmd_line
  - pyss3.util

# PySS3 - SS3 class

```python
from pyss3 import SS3

clf = SS3()

clf.fit(x_train, y_train)
y_pred = clf.predict(x_test)
```

# PySS3 - "Live Test" tool

```python
from pyss3.server import Server
from pyss3 import SS3


clf = SS3()
clf.fit(x_train, y_train)


Server.serve(clf, x_test, y_test) # <- this one! cool uh? :)
```

# PySS3 - "Live Test" tool

# PySS3 - Command Line
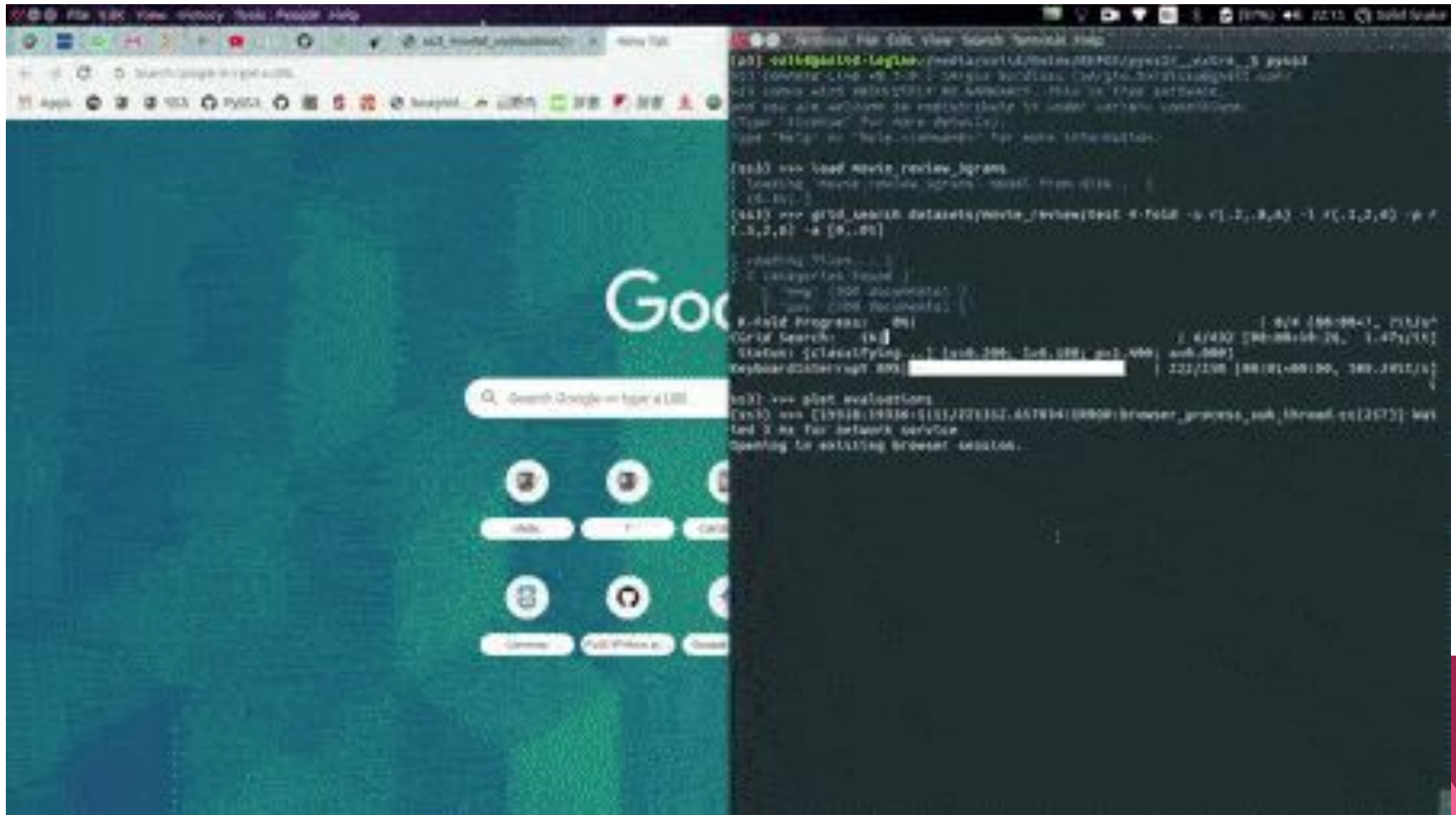
```
your@user:/your/project/path$ pyss3
(pyss3) >>> load my_model
(pyss3) >>> grid_search path/to/dataset -s [.2,.5,.8] -l [.1,1,2] -p [.5,1.5,2]
(pyss3) >>> plot evaluations
```

# PySS3 - Command Line (Evaluation Plot)

# PySS3 - ¡Manos a la obra!

- Topic Categorization

- Sentiment Analysis on Movie Reviews