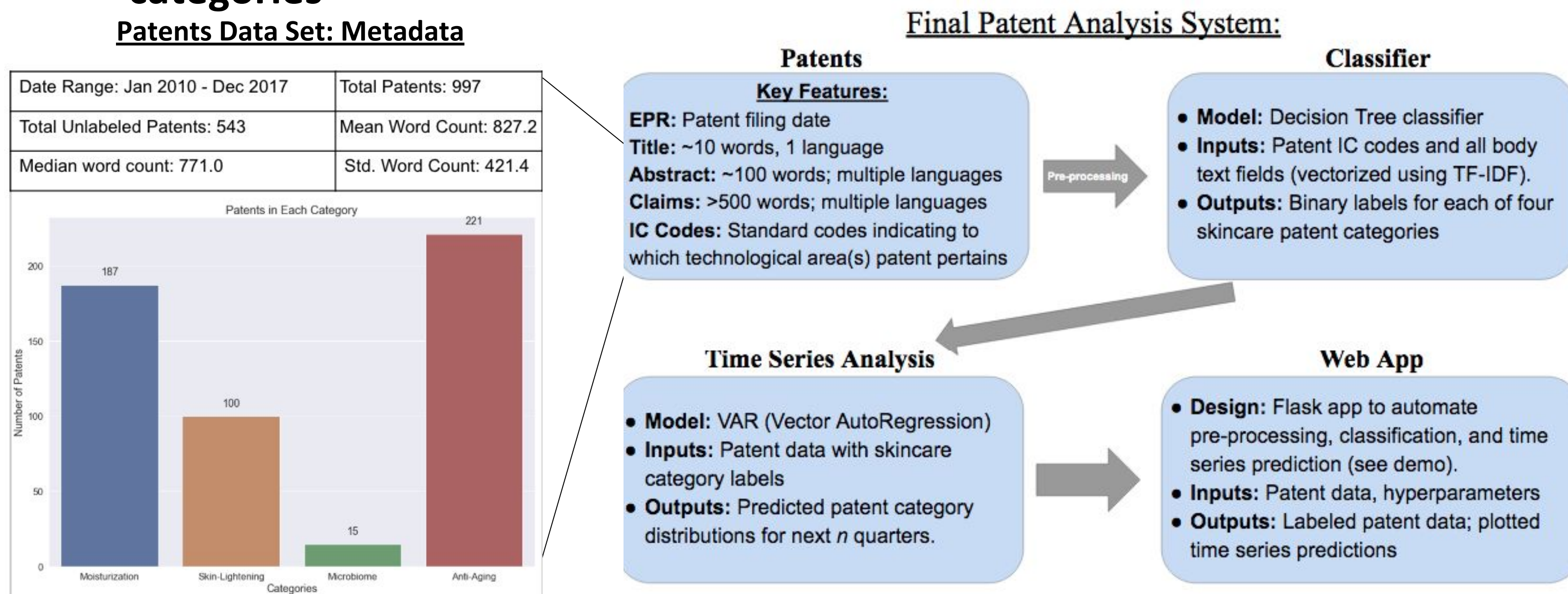


Patent Classification and Prediction System

Project Goals:

To bolster competitor surveillance capabilities, Unilever asked us to develop the following tools to analyze skincare-related patents recently filed by their biggest competitor (L'Oreal):

- A classifier to label patents as relating to any of four skincare categories: Anti-aging, Microbiome, Moisturization, and Skin Lightening
- Time series predictions for the distribution of future patents over the four categories



Patent Category Distribution - Time Series Predictions:

- Labeled data is aggregated by quarter, creating 32 data points from 2010 - 2017.
- We fit several time series models to this data and compared AIC scores.
- A VAR multivariate model, which captures interactions between categories, achieved the best overall AIC score when forecasting all categories together.
- An ARIMA model, which captures non-seasonal patterns, was best for predicting individual category distributions.

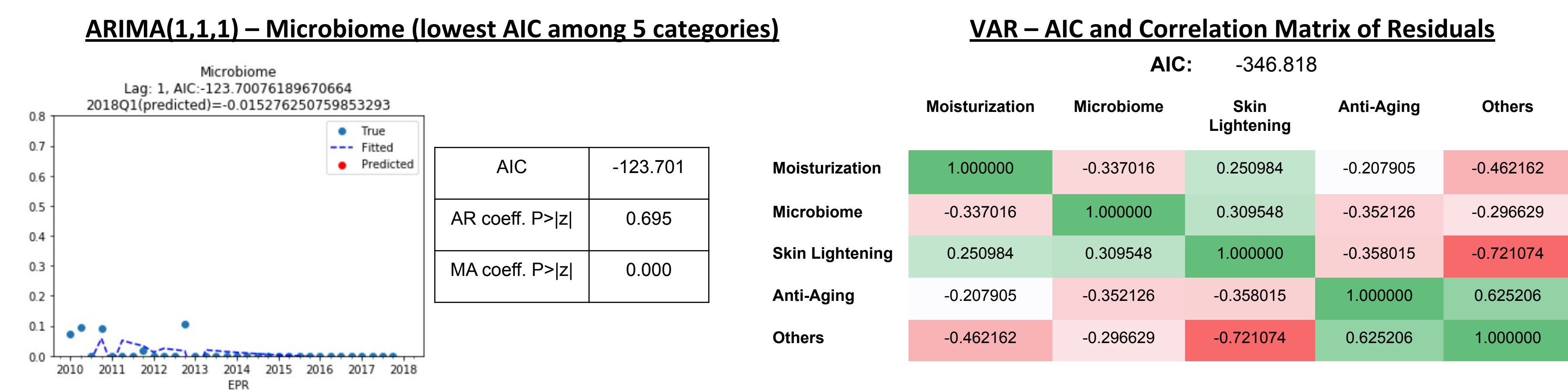


Figure 2. (Left) ARIMA fit, AIC score, and significance of AR, MA terms. (Right) VAR model AIC score and correlation matrix between category distributions.

5 Most Common Words by Skincare Category

Most common words / # of occurrences	1st	2nd	3rd	4th	5th
All Patents	composition / 1816	comprising / 1307	cosmetic / 1120	relates / 940	skin / 765
Moisturization Patents	linking / 9	reaction / 8	resistant / 7	bonds / 6	layered / 6
Microbiome Patents	container / 6	longum / 6	bifidobacterium / 4	pro / 4	alterations / 4
Skin-lightening Patents	hyperpigmentation / 4	activator / 4	eo / 3	cb / 3	receptor / 3
Anti-Aging Patents	bundle / 9	trans / 8	filters / 8	threads / 8	photoprotective / 6

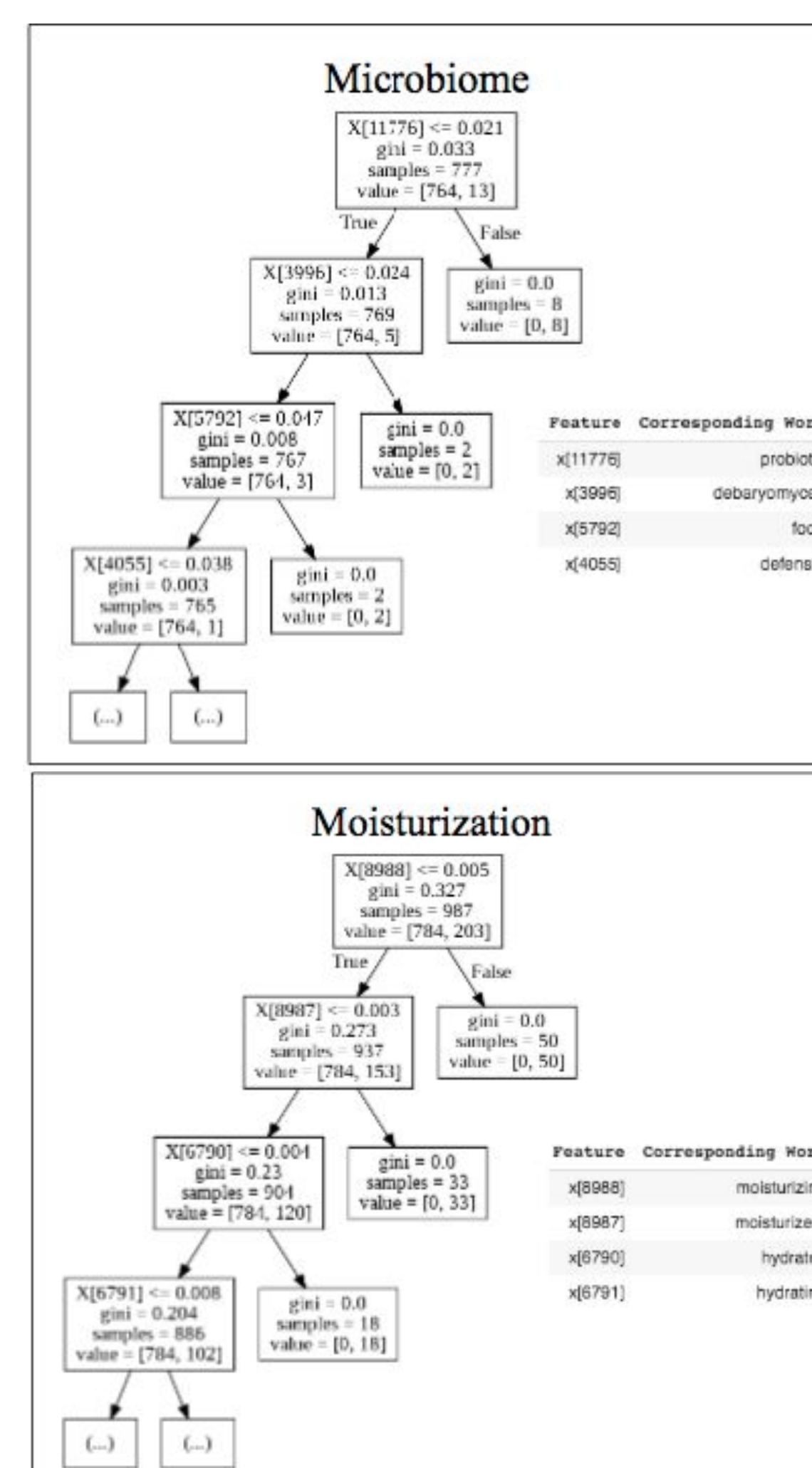
Patent Classifier:

- A pre-processing step cleans patent text, extracts its English translation only, and converts it to a TF-IDF vector.
- A patent's membership in the 'Anti-aging' and 'Skin Lightening' categories can be determined entirely from patent IC Codes, so no ML classifier is necessary.
- For 'Microbiome' and 'Moisturization' categories, we trained several models and compared their 5-fold cross-validated accuracies:

Training validation results

Model	Cross-Validated Accuracy ('Moisturization')	Cross-Validated Accuracy ('Microbiome')
Naïve Bayes (Count Vectorizer)	74.52%	96.63%
Naïve Bayes (TF-IDF Vectorizer)	79.44%	98.50%
Decision Tree	90.27%	98.51%
Random Forest	81.14%	--
SVM	79.44%	--

Figure 1. Decision tree structure visualization



Web Application:

- The user uploads a csv containing conforming patent data, as specified on the app's home page.
- The app classifies the data and returns a new csv with binary category indicator fields appended.
- The user can check this output, then upload it again to run time series analysis and obtain predictions as a plot and a table.
- Check out our demo application!

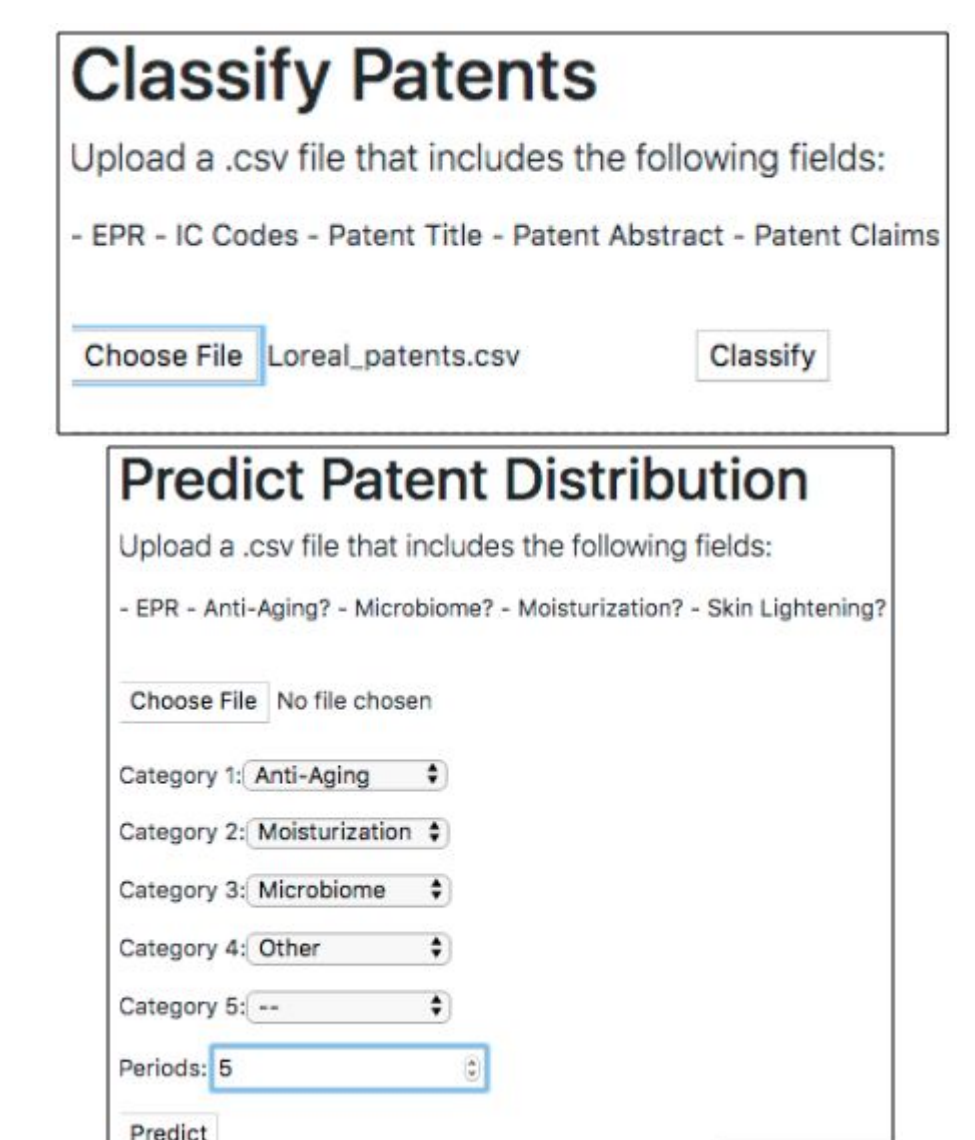


Figure 3. Web app homepage

Conclusion:

- Our models, accessed via the web app, offer Unilever an intelligent way to track competitor patent activity.
- With more manually labeled training data examples, our classification and time series prediction models could potentially be improved.
- We will continue working with Unilever to develop additional visualization and analytical capabilities of the web app.

Acknowledgments:

Special thanks to Kriste Krstovski and Vikash Khanna for their collaboration and support!

References:

- [1] C. Bishop, *Pattern Recognition and Machine Learning*, Springer.
- [2] Scikit-learn: Machine Learning in Python, Pedregosa *et al.*, *JMLR* 12, pp. 2825-2830, 2011.