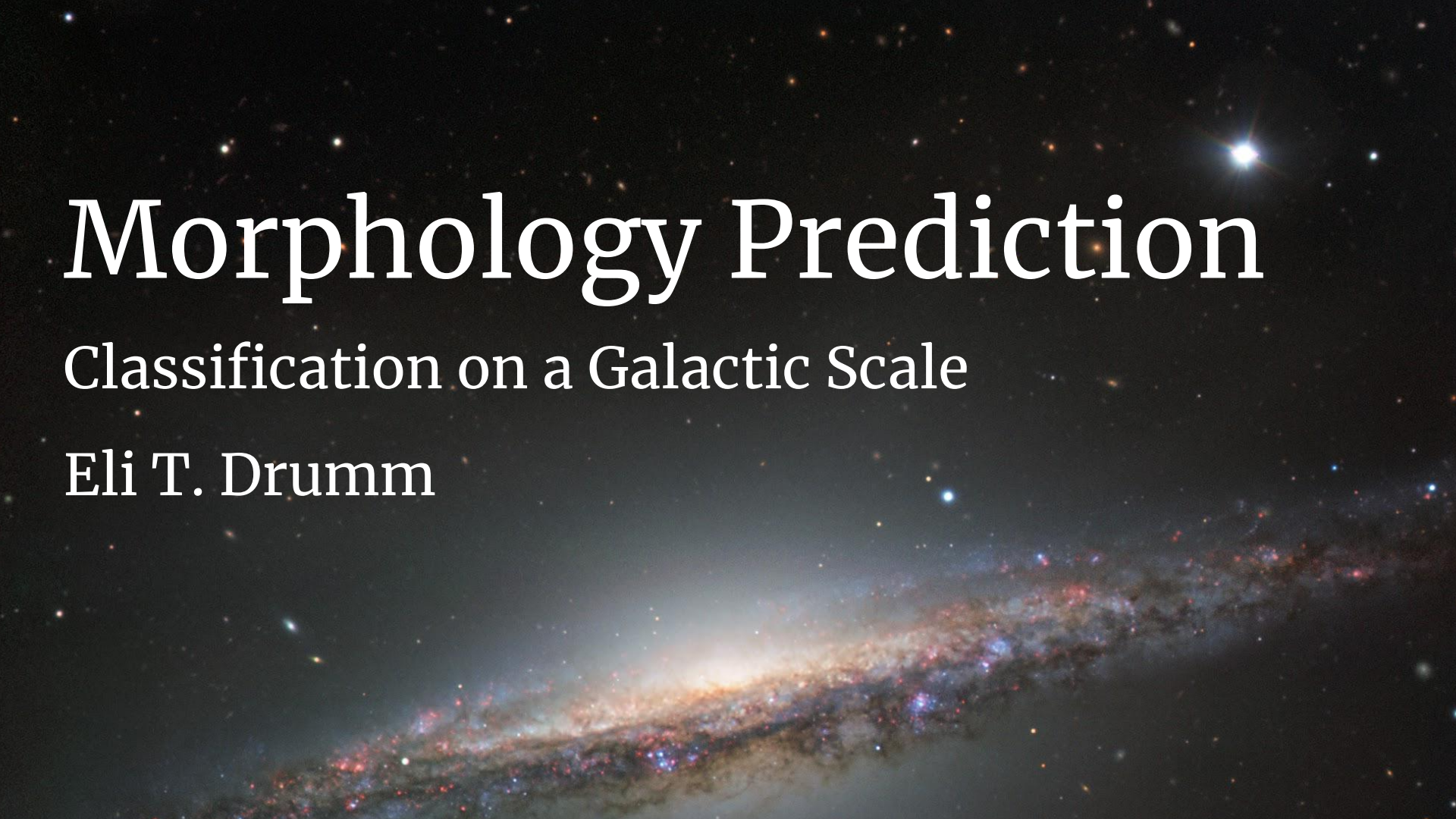


Morphology Prediction

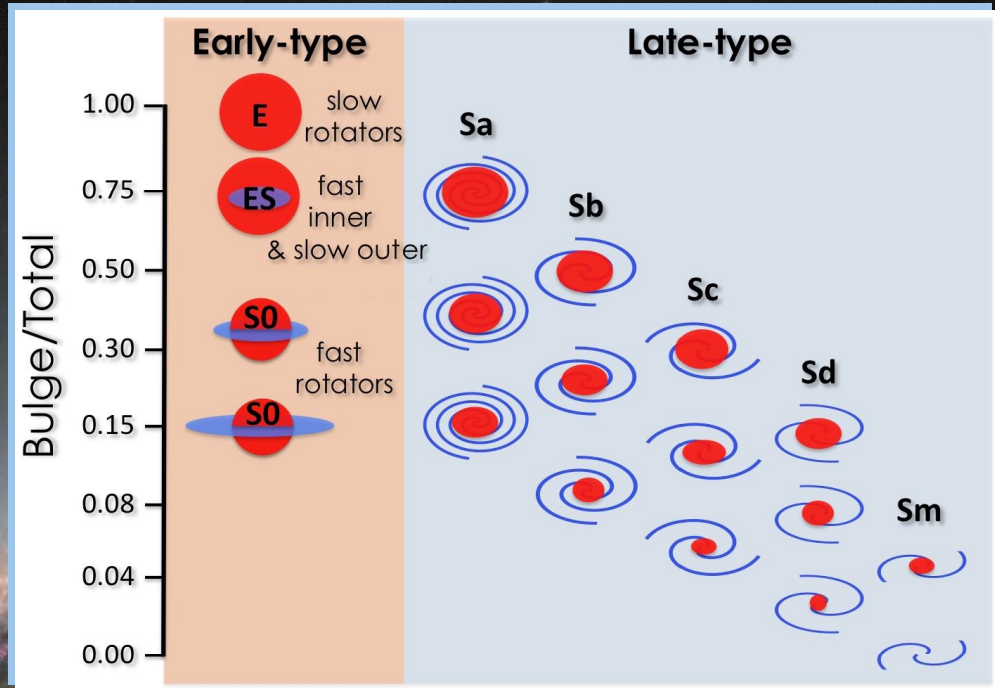
Classification on a Galactic Scale

Eli T. Drumm



Background

- Galaxy Zoo project
- Paper (and dataset) by Schawinski et al. (2009)
- Early- vs. late-type



Feature selection

- Spectroscopy
 - *ugriz*
- M_{\odot}
- Redshift
- Luminosity
- etc.



Linear model

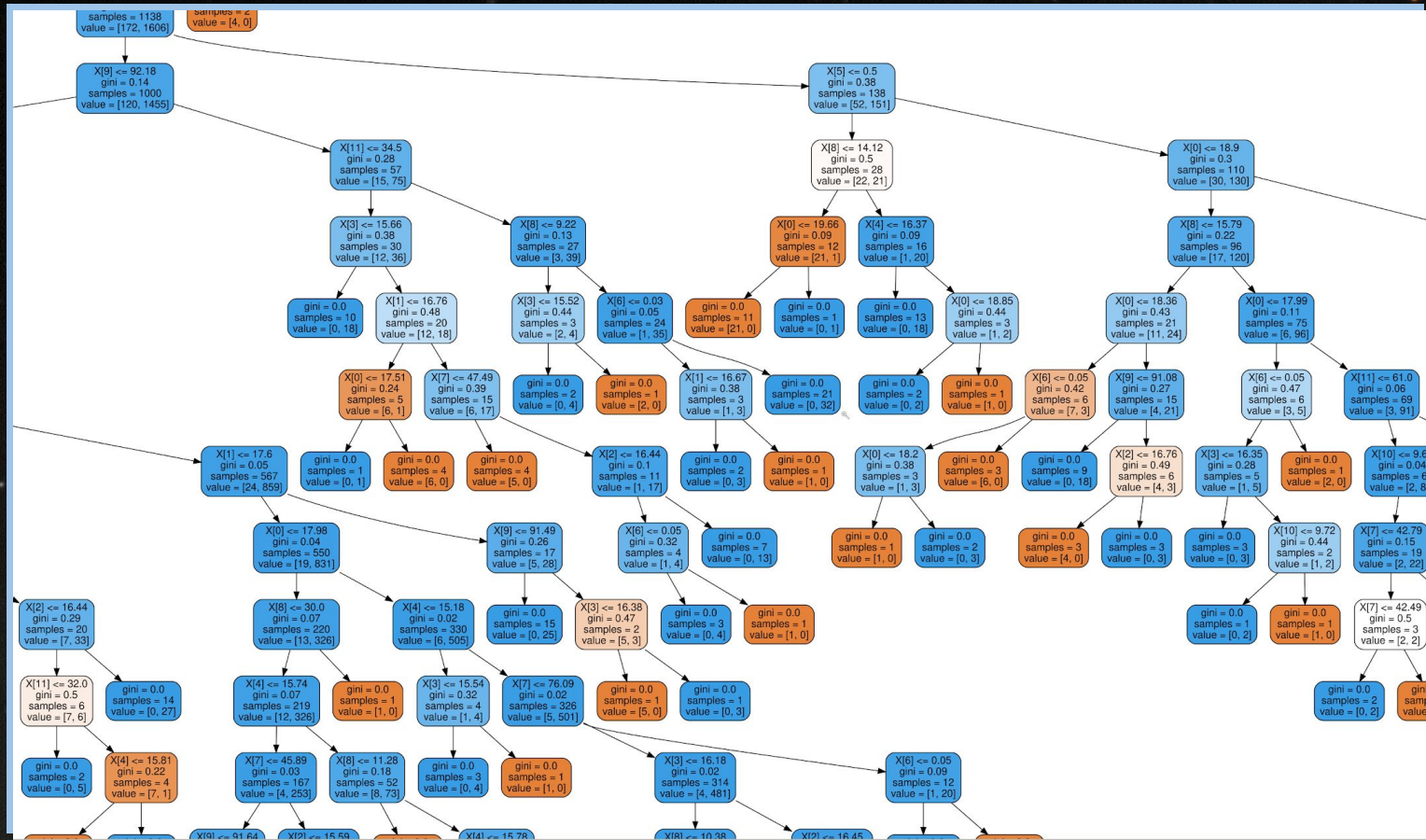
- Logistic regression, 5-fold CV
 - aaaand...
- Accuracy is 87%
- ROC AUC is 0.923

	precision	recall	f1
non-late	0.85	0.78	0.81
late-type	0.88	0.92	0.90

Trees and forests

- Can we do better?
- Random forest
 - Again, with 5-fold CV





Random forest results

Accuracy is 89.6%.

ROC AUC 0.942.

	precision	recall	f1
non-late	0.87	0.82	0.85
late-type	0.90	0.92	0.91

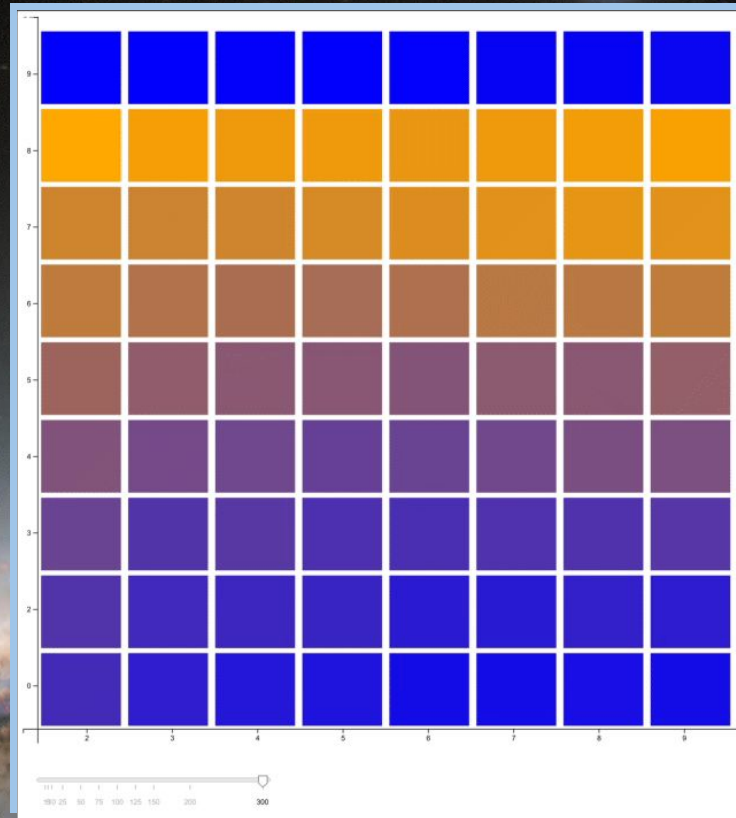


Hiking the random forests

We can do a grid search to find the best parameters.

General lessons:

- Plant more trees
- ...but expect diminishing returns
- Unlimited max depth

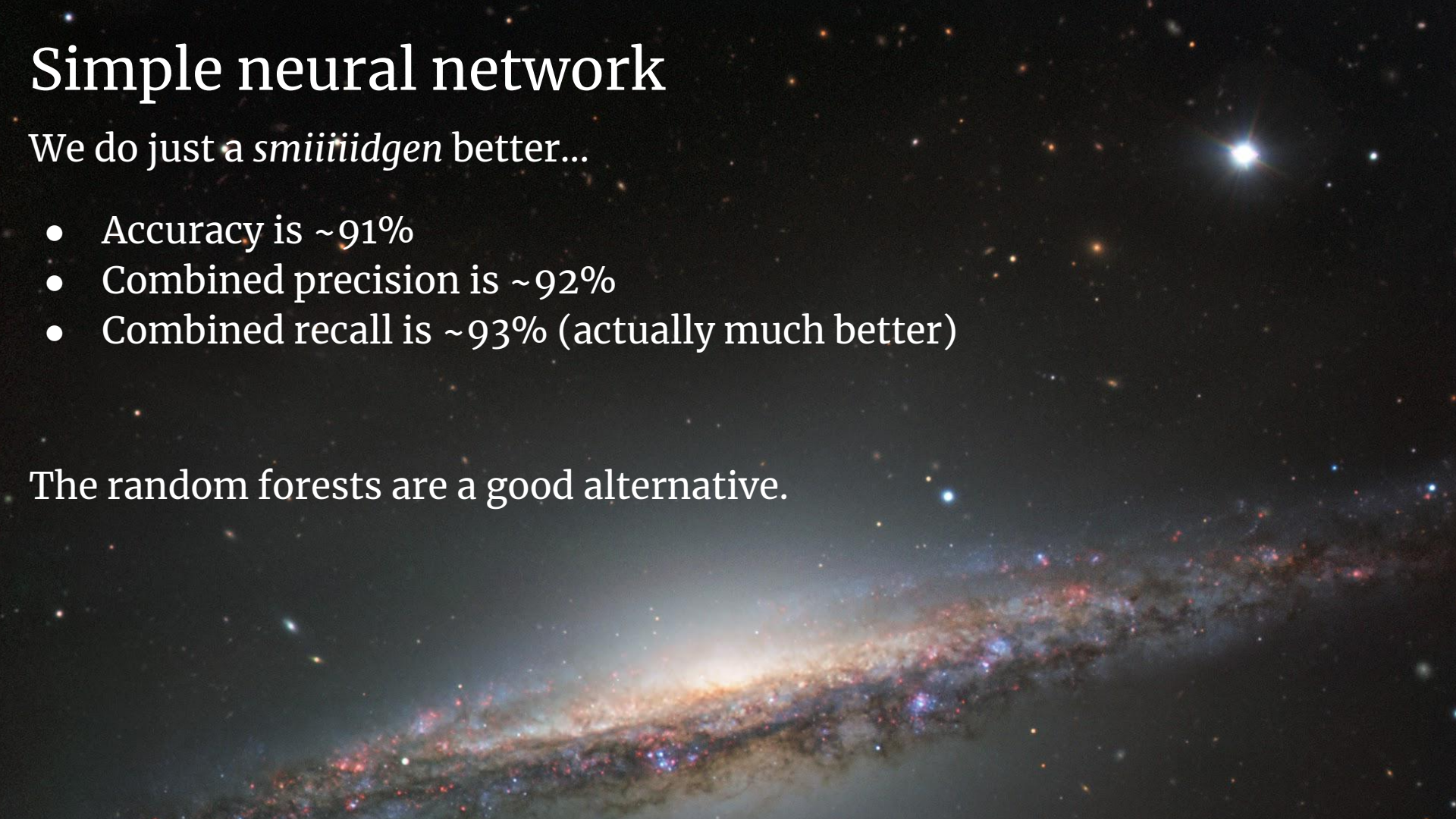


Simple neural network

We do just a *smiiiidgen* better...

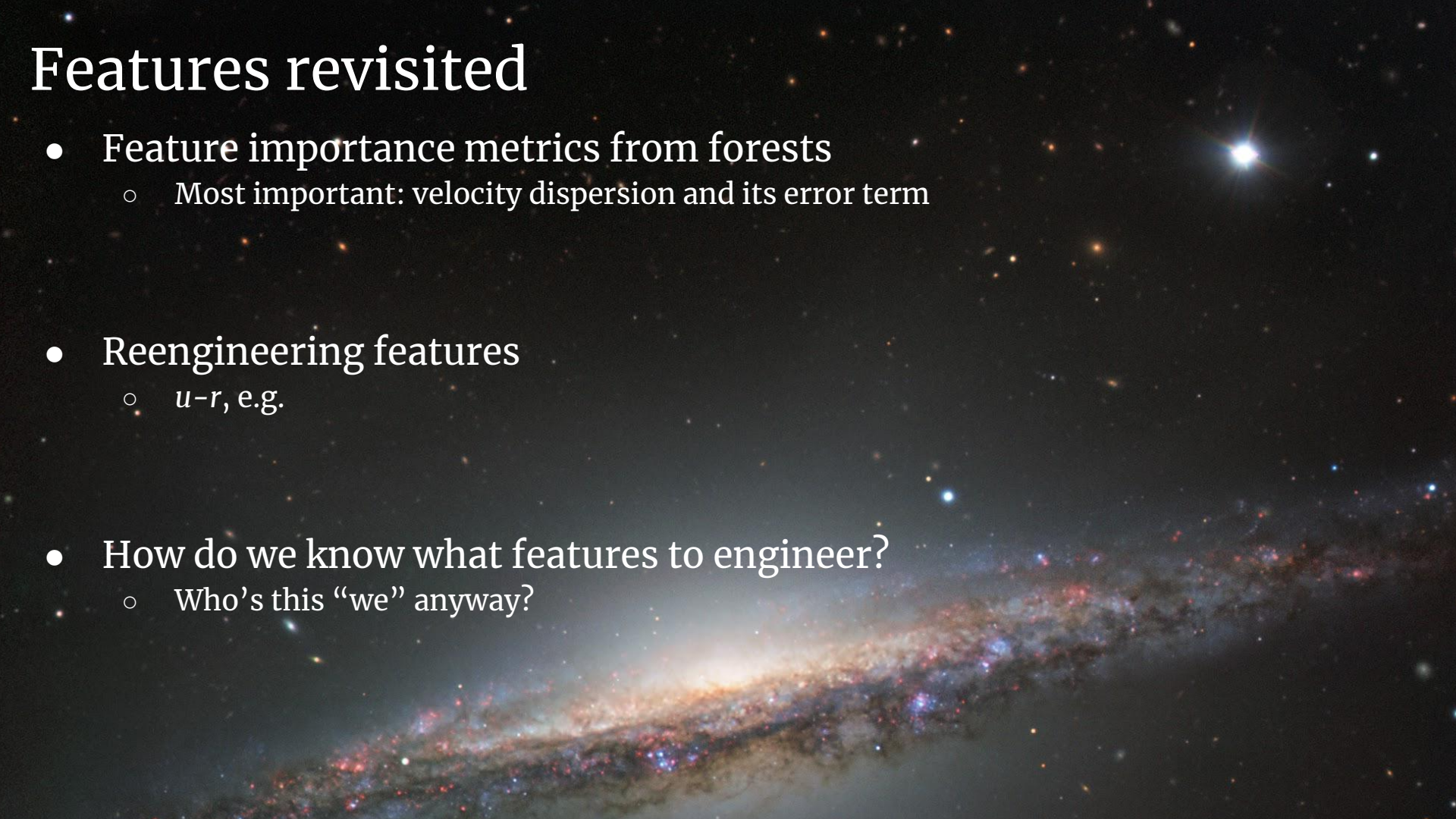
- Accuracy is ~91%
- Combined precision is ~92%
- Combined recall is ~93% (actually much better)

The random forests are a good alternative.



Features revisited

- Feature importance metrics from forests
 - Most important: velocity dispersion and its error term
- Reengineering features
 - $u-r$, e.g.
- How do we know what features to engineer?
 - Who's this “we” anyway?



Takeaways

- Other data: Galaxy Zoo 2, Galaxy Zoo Hubble
- Without computed data, prediction is difficult
 - *“At the best galaxy redshift, the stellar velocity dispersion is also determined. This is done by computing a PCA basis of 24 eigenspectra from the ELODIE stellar library (Prugniel & Soubiran 2001), convolved and binned to match the instrumental resolution and constant-velocity pixel scale of the reduced SDSS spectra, and broadened by Gaussian kernels of successively larger velocity width ranging from 100 to 850 km/s in steps of 25 km/s. The broadened stellar template sets are redshifted to the best-fit galaxy redshift, and the spectrum is modeled as a least-squares linear combination of the basis at each trial broadening, masking pixels at the position of common emission lines in the galaxy-redshift rest frame. The best-fit velocity dispersion is determined by fitting locally for the position of the minimum of chi-squared versus trial velocity dispersion in the neighborhood of the lowest gridded chi-squared value. Velocity-dispersion error estimates are determined from the curvature of the chi-squared curve at the global minimum, and are set to a negative value if the best value occurs at the high-velocity end of the fitting range. Reported best-fit velocity-dispersion values less than about 100 km/s are below the resolution limit of the SDSS spectrograph and are to be regarded with caution.”*
- Morale: data science can't be isolated
 - Work with astronomers or whomever!