

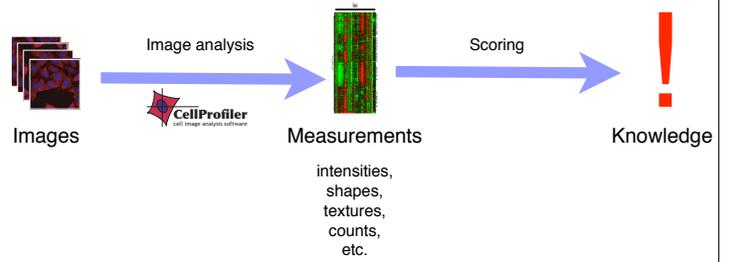
Large-scale learning of cellular phenotypes from images

Vejbjorn Ljosa, Piyush B. Gupta, Thouis R. Jones,
Eric S. Lander, and Anne E. Carpenter

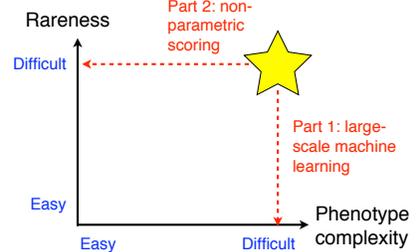
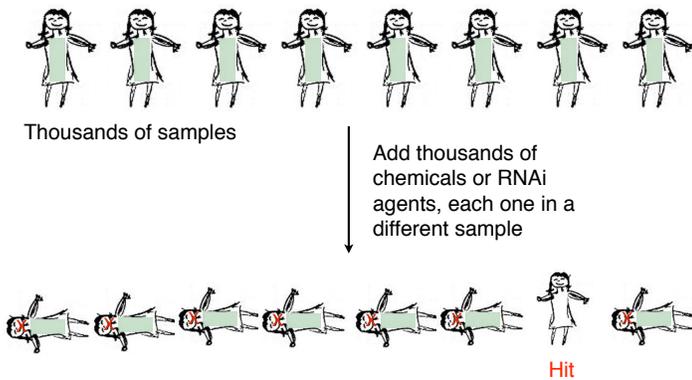


RECOMB Systems Biology
Cambridge, MA – 2009-12-06

From images to knowledge in high throughput



Screens

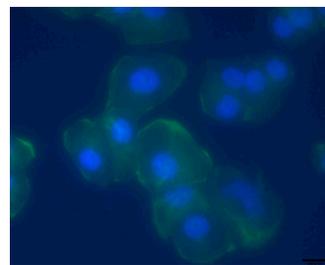


- Simple methods
- Few assumptions
- Little modeling
- ⋮
- ⋮
- Lots of data

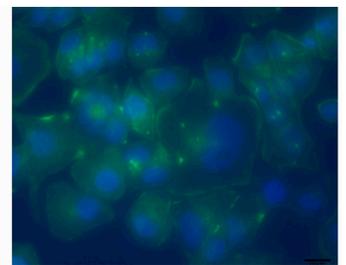
The phenotype of motile T47D cells

Normal T47D cells

Features associated with cell motility: lamellipodia, filopodia, polarized cell shape, F-actin nucleation at filopodia, less clumping

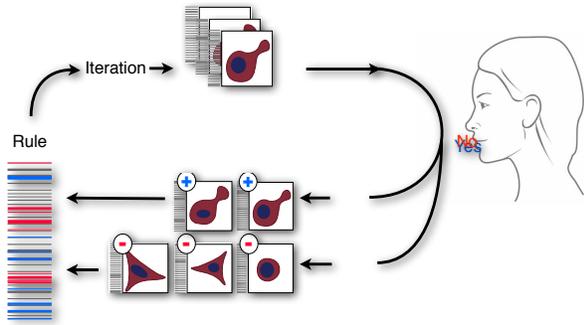


Unstimulated



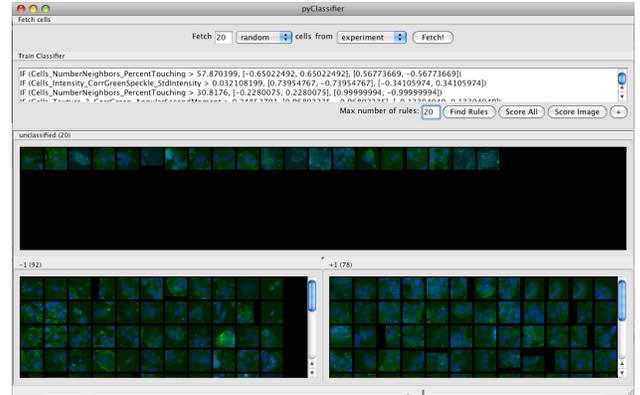
Stimulated by heregulin

Iterative machine learning

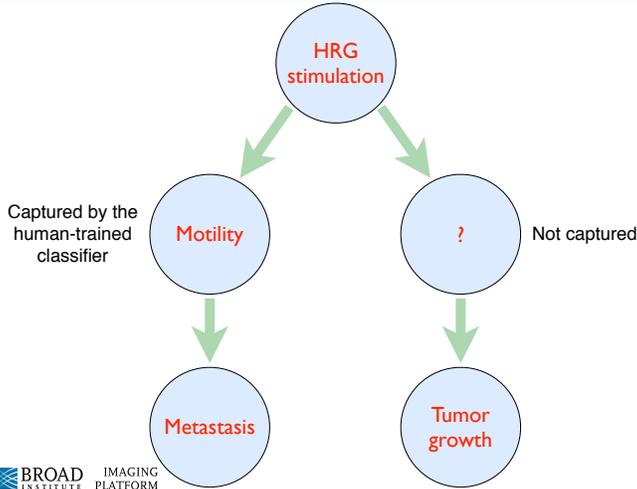


CellProfiler Analyst [Jones et al., PNAS, 2009]
data exploration software
Using gentle boosting [Friedman et al., 1998]

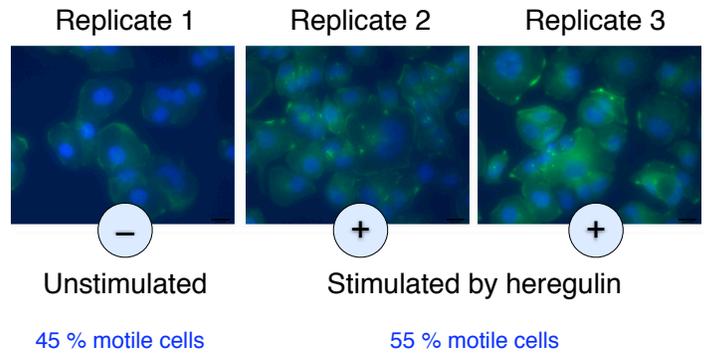
Built training set of ~300 cells



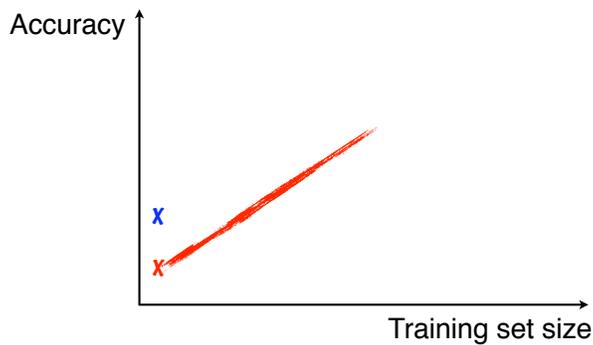
Why cut out the human?



Labeling for automatic training set

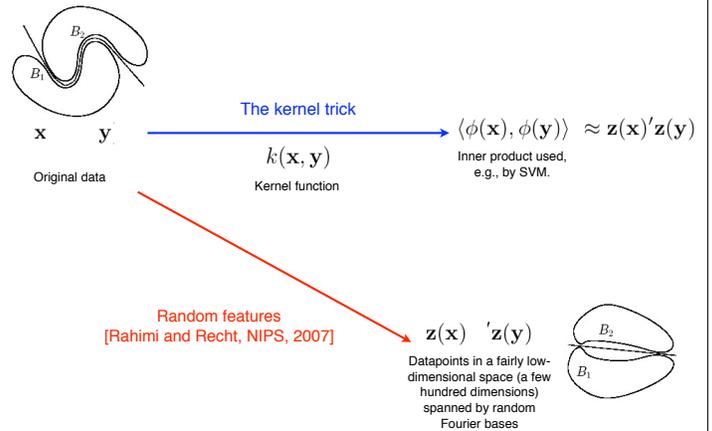


Two ways to improve the classifier



See [Banko & Brill, 2001]

Random features

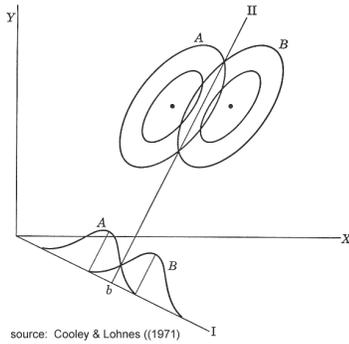


Linear discriminant on random features

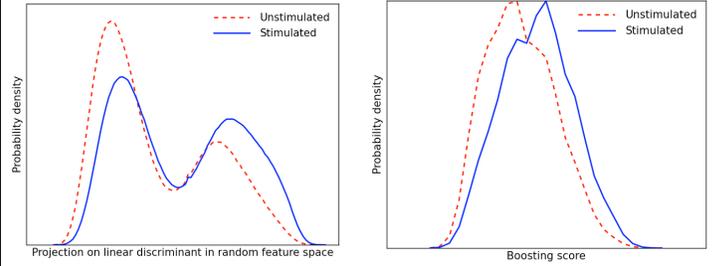
7.6 million training cells,
130 measurements

Mapped into 250-
dimensional random
feature space

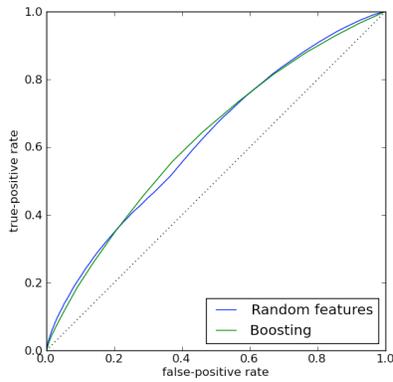
Trained Fisher's linear
discriminant



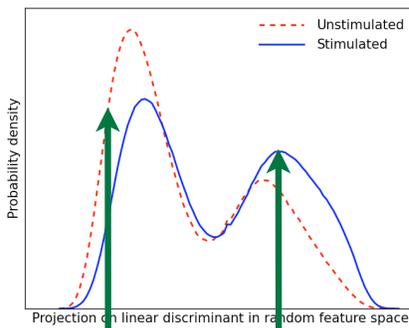
Automatic vs. hand training



Automated classifier as good as human-trained



Soft labels for cells



Cell with this score: $P(\text{stimulated}) = .85$
Cell with this score: $P(\text{stimulated}) = .65$

Fuzzy counts

Each cell in a sample is assigned a
probability of having the positive phenotype



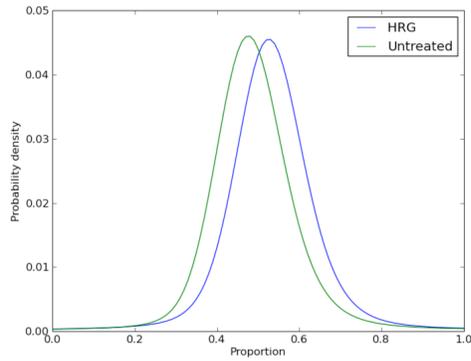
Compute the probability of the sample having

- 0 positive cells
- 1 positive cell
- 2 positive cells
- ...
- 6 positive cells

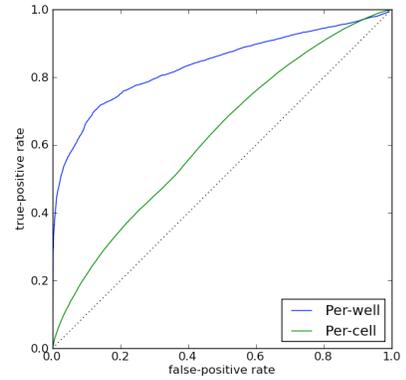
This pdf of counts can be turned into a pdf of proportions

Scoring samples by fuzzy counts

Mix the pdfs of proportions => empirical positive and negative control distributions

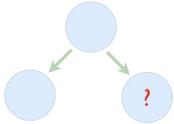
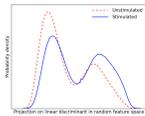


Per-well accuracy is quite good



Summary

We can identify subtle, complex cellular phenotypes without human training



May enable screening for “invisible” phenotypes, as well as large-scale profiling experiments

Avoid premature thresholding, classification, and aggregation. Embrace populations, uncertain values, and fuzzy scores.

