# From Zero to Hero

### Anticipating Zero Results From Query Features, Ignoring Content

*Mikhail Popov (Analysis & Report)*

*May 17, 2016*

### Executive Summary

The Discovery Department uses the zero results rate – the proportion of searches that yield zero results – to measure the performance of our search system. However, little is known about possible patterns that affect the quantity (and quality) of results our users see. In this report, we use random forest and logistic regression models to shed light on the types of queries that tend to yield zero results.

Namely, we found that whether the query has an even number of double quotes is one of the most important indicators of whether it will yield zero results. Other notable features that impact the quantity of results include: whether the query is only punctuation and spaces, whether it ends with ?, and whether it has logical operators. For a full list of features and their importance and impact, please see Figures 4 and 5.

## Introduction

One of Discovery/Search team's key performance indicator (KPI) metrics is the zero results rate – the proportion of searches that yield zero results. However, little is known about possible patterns that affect the search engine as it retrieves potential matches for the user-provided query. If we don't know which types of queries are more or less likely than others to yield zero results, that makes it extremely challenging to come up with configuration tweaks and other modifications that address particular types of queries, and ultimately help the user discover the content they are searching for.

This report attempts to bring to light those patterns which affect the chances of getting zero or some results. First we deconstruct a large set of queries into descriptive features. Then we use two classification methods which allow us to make inference on those extracted features. Random forests enable us to assess how important certain features are (through various variable importance measures) in classification, while logistic regression enables us to asses the magnitude and direction of a feature's impact on the probability of zero results.

## Methods

On 24 February 2016 we completed a Hive user-defined function (UDF) for deconstructing search queries into various features (see Patch 254461). The UDF detects a variety of features such as: odd or even number of quotation marks, logical operators (e.g. AND, OR), "prefix:" or "insource:", and wildcards. For a full list of features, please see T118218 and SearchQuery.java and SearchQueryFeatureRegex.java source code.

We extracted 5.7 million queries made on the web – and excluding known automata – to train a random forest (an ensemble classification algorithm) on queries yielding zero or some results. We decided to grow a random forest because this particular classification algorithm allows us to assess

the importance of different features (see variable importance). We deconstructed the queries into "query features" (yielding an indicator matrix of 1s and 0s) and combined that with other statistics (number of terms, number of characters, and number of features detected) to construct a final feature matrix used to grow the random forest. We considered including country as a predictor but early parameter tuning tests showed that a small increase in prediction accuracy was not enough to offset the time it would take to train the algorithm on a feature space with 234 additional dimensions (one per each country observed).

A random 80% subset of the data was used to grow the trees, with the remainder 20% as out-of-bag (OOB) data for assessing classification error and variable importance. In particular, we use the *Mean Decrease Gini* (MDI) and *Mean Decrease Accuracy* (MDA) to find the features which have the most impact on classification, and thus may provide us with ideas about what we can tweak to make search better. MDA works such that if a variable $X_j$ is associated to response $Y$, then randomly permuting its values should result in a substantial decrease in the accuracy of the classification.

Additionally, we used logistic regression on 1% of the data to assess the magnitude and the direction of the features' effects on the query yielding some or zero results.

## Results

Figures 1–3 on the following pages break down the zero results by country, extracted feature, and combination of extracted features. Figures 4 and 5 in the Random Forest and Logistic Regression sections, respectively, compare the impact those features have in the context of predicting zero results. Figure 6 compares the features using the measures from both models to reveal which features the two models agree on.

## References

Random forest. URL **https://en.wikipedia.org/wiki/Random_forest**

Louppe, G. (2014, July 28). Understanding Random Forests: From Theory to Practice.

Hastie, T., Tibshirani, R., & Friedman, J. (2009). The Elements of Statistical Learning. New York, NY: Springer Science & Business Media. **http://doi.org/10.1007/978-0-387-84858-7**

A. Liaw and M. Wiener (2002). Classification and Regression by randomForest. R News 2(3), 18–22.

R Core Team (2016). R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. URL **https://www.R-project.org/**.

## Zero results rate in top 20 countries by volume of searches
Search queries made on the web, excluding known automata on Sun (01 May 2016)

| Country | |
|---------|---|
| Ireland | 11.7% of 312.7K |
| Italy | 13.1% of 139.1K |
| United Kingdom | 15.2% of 185.9K |
| Angola | 15.9% of 75.5K |
| Brazil | 16.9% of 94.4K |
| France | 17.4% of 420.3K |
| United States | 17.6% of 1.2M |
| Germany | 17.9% of 437.1K |
| India | 19.8% of 160K |
| Poland | 19.9% of 78.4K |
| Czech Republic | 20.3% of 68.4K |
| Russia | 21.1% of 139.9K |
| Netherlands | 21.8% of 84.4K |
| Iran | 22.4% of 66.7K |
| Australia | 23.8% of 71.3K |
| Spain | 27.1% of 108.9K |
| Japan | 29.7% of 157.8K |
| South Korea | 30.6% of 67.3K |
| China | 30.9% of 65.3K |
| Canada | 34.2% of 133.2K |

Countries (y-axis) · Zero Results Rate (x-axis: 0%, 5%, 10%, 15%, 20%, 25%, 30%, 35%, 40%)

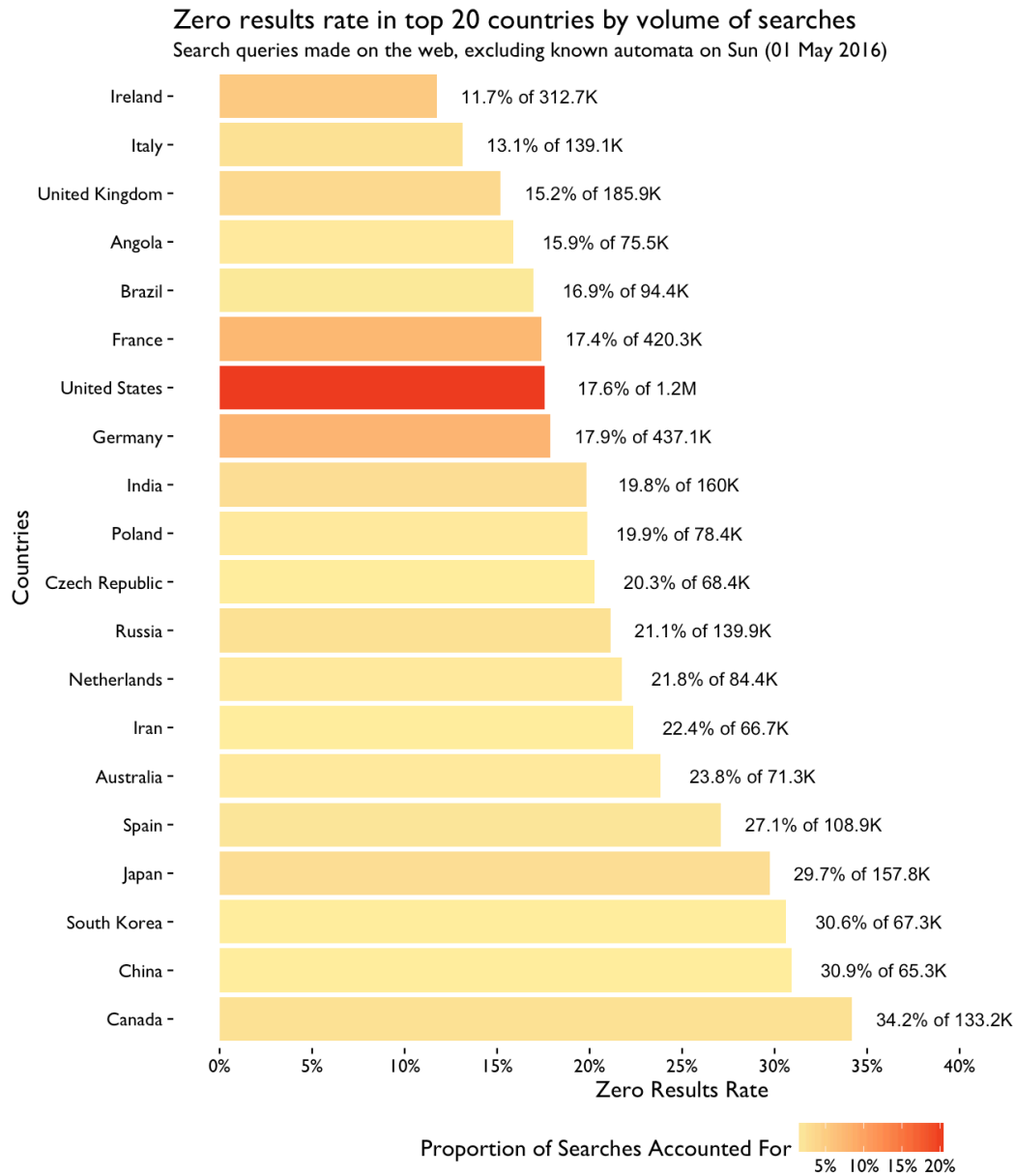Proportion of Searches Accounted For — 5% 10% 15% 20%

*Figure 1: Zero results rate by country, limited to top 20 countries by volume of searches performed. More than 20% of the searches were from United States, which had a ZRR of 17.6%. Canada had the highest zero results rate but only accounted for 2.3% of the searches in the dataset.*

3

## Proportion of searches with zero results by query feature
Search queries made on the web, excluding known automata

| Feature | Zero Results Rate |
|---|---|
| is just ~ | 0.0% of 1 |
| is empty | 0.0% of 88 |
| has OR | 14.7% of 29.6K |
| is simple | 16.2% of 4.4M |
| is incategory | 16.2% of 2K |
| is intitle | 21.9% of 5.6K |
| has logic inversion (-) | 23.7% of 3.6K |
| is just ? | 25.0% of 16 |
| has non-ASCII | 25.4% of 1.2M |
| has NOT | 28.0% of 25 |
| is fuzzy search | 32.4% of 139 |
| is insource | 33.2% of 653 |
| is prefix | 33.8% of 3.4K |
| has logic inversion (!) | 34.1% of 179 |
| is just * | 37.5% of 8 |
| has AND | 42.8% of 1.7K |
| has odd double quotes | 43.7% of 3.7K |
| has one double quote | 45.3% of 3.3K |
| has wildcard | 52.3% of 34.5K |
| ends with ? | 56.6% of 18K |
| has quot | 61.5% of 13 |
| has even double quotes | 76.1% of 210.9K |
| is only punctuation and spaces | 83.4% of 2.9K |
| has @ | 84.5% of 4.1K |
| forces search results | 100.0% of 1 |

Number of Queries: 10 · 100 · 1K · 10K · 100K · 1M

*Figure 2: Zero results rate by extracted feature, taken individually.*

4

## Proportion of searches with zero results by query features (sorted by ZRR)
Search queries made on the web, excluding known automata

| Features | Zero Results Rate |
|---|---|
| is incategory, has logic inversion (-), has even double quotes | 0.0% of 106 |
| has OR, has wildcard, has non-ASCII | 0.0% of 439 |
| has OR, has odd double quotes | 0.0% of 218 |
| has OR, has wildcard | 2.1% of 2.1K |
| is intitle, has even double quotes | 2.5% of 1.7K |
| has OR | 5.3% of 4.4K |
| is incategory, is intitle, has even double quotes | 9.7% of 176 |
| has OR, has non-ASCII | 11.3% of 851 |
| has OR, has non-ASCII, has even double quotes | 12.0% of 4.7K |
| has OR, has even double quotes | 12.9% of 12.6K |
| is incategory, has non-ASCII, has even double quotes | 14.7% of 577 |
| is incategory, has even double quotes | 15.3% of 1K |
| is simple | 16.2% of 4.4M |
| has logic inversion (-), has non-ASCII | 18.9% of 783 |
| has logic inversion (-) | 21.1% of 2.1K |
| has logic inversion (!) | 23.7% of 131 |
| has non-ASCII | 24.6% of 1.1M |
| is intitle | 28.2% of 3.4K |
| is prefix | 31.6% of 2.5K |
| has AND | 33.7% of 1.4K |
| has one double quote, has odd double quotes | 38.9% of 2.1K |
| has OR, has wildcard, has non-ASCII, has even double quotes | 40.0% of 1K |
| is prefix, has non-ASCII | 41.0% of 790 |
| has wildcard, has non-ASCII, has even double quotes | 41.4% of 198 |
| is intitle, has non-ASCII | 44.2% of 292 |
| has OR, has wildcard, has even double quotes | 44.8% of 2.9K |
| is insource | 45.0% of 436 |
| has logic inversion (-), has even double quotes | 46.9% of 369 |
| has non-ASCII, has one double quote, has odd double quotes | 47.7% of 704 |
| has wildcard, has even double quotes | 47.8% of 801 |
| ends with ?, has wildcard | 52.2% of 12.7K |
| has logic inversion (-), has non-ASCII, has even double quotes | 56.0% of 109 |
| ends with ?, has wildcard, has even double quotes | 56.5% of 115 |
| has wildcard | 60.9% of 7.4K |
| ends with ?, has wildcard, has non-ASCII | 64.9% of 4.5K |
| has wildcard, has non-ASCII | 73.5% of 1.6K |
| is only punctuation and spaces | 78.2% of 2.1K |
| has non-ASCII, has even double quotes | 78.8% of 13.1K |
| has @ | 82.1% of 2.1K |
| has @, has non-ASCII | 83.8% of 321 |
| has even double quotes | 84.8% of 170.3K |
| has @, has wildcard | 86.2% of 116 |
| has @, has even double quotes | 97.4% of 392 |
| ends with ? | 97.8% of 459 |
| has @, is only punctuation and spaces | 98.6% of 693 |

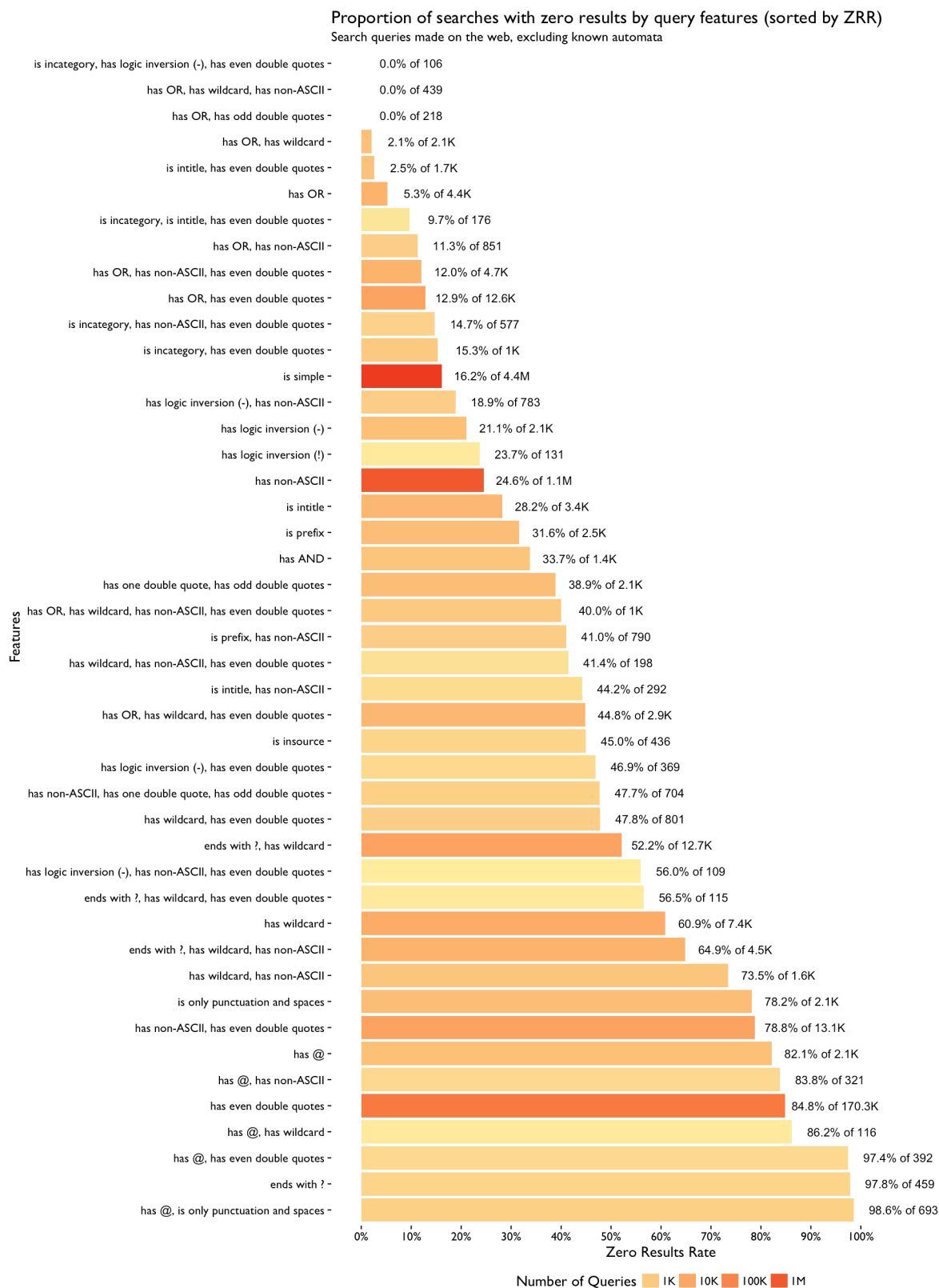Number of Queries: 1K, 10K, 100K, 1M

*Figure 3: Zero results rate by combinations of extracted features. For space and reliability reasons, we restricted this chart to only show combinations which had more than 100 queries. A vast majority of the combinations include the "has even double" feature.*

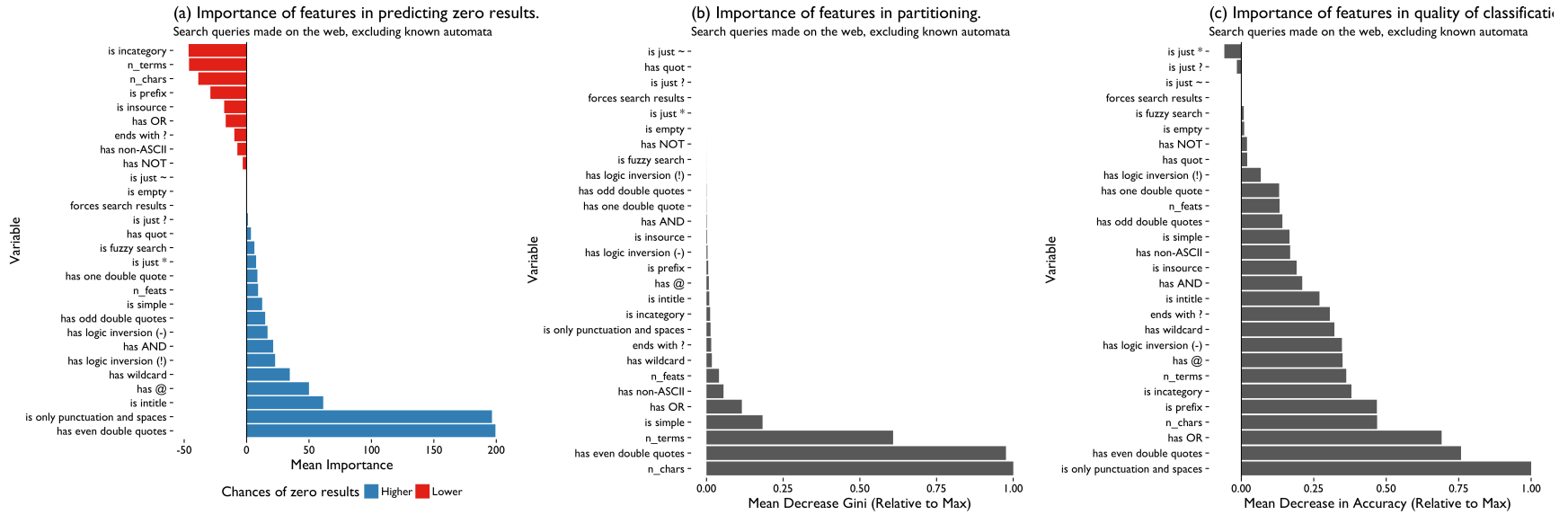# Classification via Random Forest



*Figure 4: **(a)** Variable importance according to mean decrease in accuracy (increase in prediction error after permuting values of the predictor) specific to zero results queries. **(b)** Variable importance according to mean decrease in impurity, using Gini index. **(c)** Variable importance according to mean decrease in accuracy over both classes of queries (zero results and some results).*

In Figure 4 above, we can see that across the battery of variable importance metrics, having an even number of double quotes is 1 of the 2 most important indicators of whether the query will have some or zero results. This makes sense because the engine will look for exact matches to each quoted segment of the query, meaning less documents will be retrieved. Other high-importance features make a lot of logical sense, such as whether the query includes the OR operator or whether it's only punctuation and spaces.

**Classification via Logistic Regression**



*Figure 5: Point estimates and 95% confidence intervals for coefficients used in the logistic regression model. Positive coefficients increase the odds of the query yielding zero results.*

The two models reveal a similar set of high-impact features, namely:

- Queries with an even number of " marks
- Queries that are just punctuation and spaces (obvious)
- Queries with wildcards (less obvious)
- Queries that end with ?
- Queries with @ (emails, possibly?)
- Queries that have the operators AND and OR

In Figure 6 below, we can see which features the two models agree on as having a relationship with the results the query yields. Specifically, they agree on the following features that are important – with respect to various mean decrease accuracies (MDAs) – and make zero results more likely: has even double quotes, has the at symbol, is prefix, has logic inversion, and ends with a question mark.
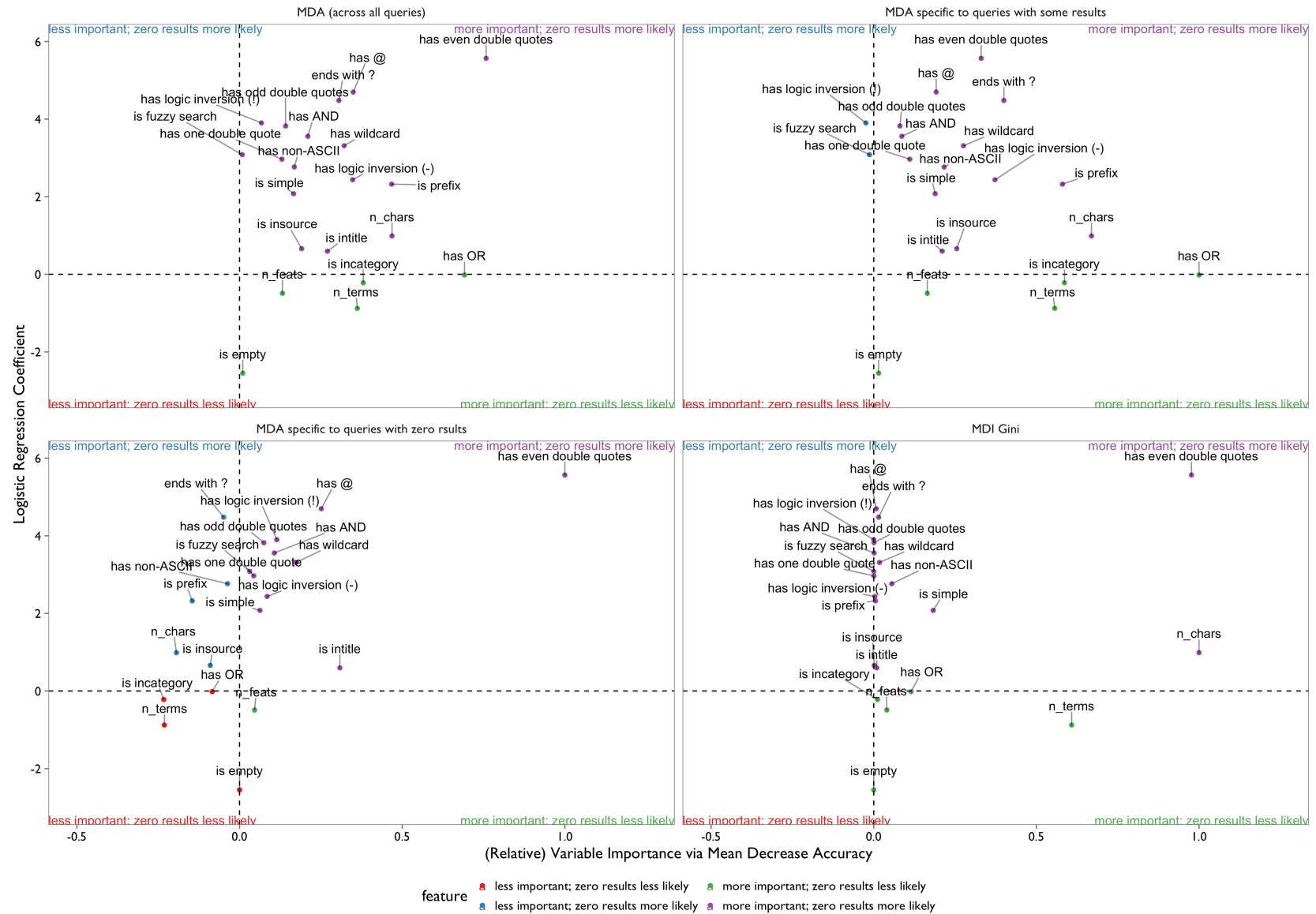
Figure 6: A scatter map of features with respect to variable importance (via relative mean decrease metrics) and logistic regression coefficient estimates. The plot is divided into quadrants according to how important or unimportant the feature is in random forest classification and whether a query having the feature is more or less likely to yield zero results.