

PyData New York City 2017

—

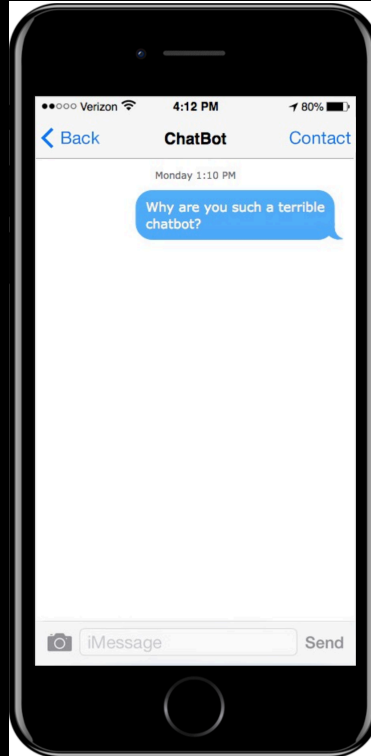
Nick Acosta
Developer Advocate

Nick Acosta

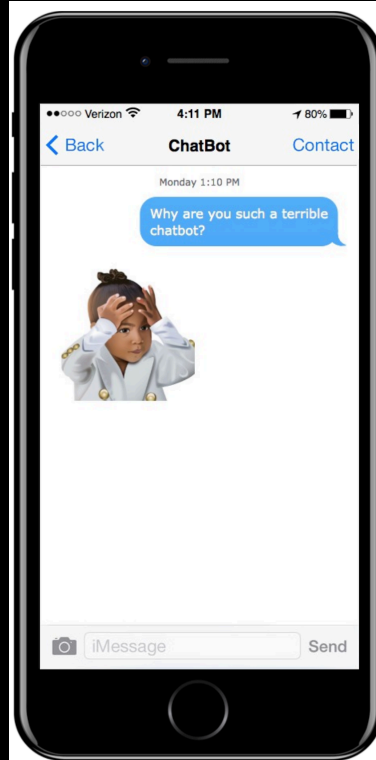
Before becoming a Developer Advocate at IBM, Nick studied computer science at Purdue University and the University of Southern California, and was a high performance computing consultant for Hewlett-Packard in Grenoble, France. He now specializes in machine learning and interacting with other data scientists of various communities, startups, and enterprises in order to help them succeed on IBM's data science platform. He has a strong interest in data science education and all things Kardashian.

Motivation for talk chatbots that respond with images

Motivation for talk chatbots that respond with images



Motivation for talk chatbots that respond with images



Attachment upload API

Attachment upload API allows you to upload an attachment that you may later send out to many users, without having to repeatedly upload the same data each time it is sent :

```
var attachment = {
  "type": "image",
  "payload": {
    "url": "https://pbs.twimg.com/profile_images/803642201653858305/IAW1DBPw_400x400.png",
    "is_reusable": true
  }
};

controller.api.attachment_upload.upload(attachment, function (err, attachmentId) {
  if(err) {
    // Error
  } else {
    var image = {
      "attachment": {
        "type": "image",
        "payload": {
          "attachment_id": attachmentId
        }
      }
    };
    bot.reply(message, image);
  }
});
```

What am I
even looking
at?

Agenda

Part One - Introduction

About me	01
Topic Introduction	02
Agenda	03
	08

Part Two - The Data

Data selection	09
PyGitHub	10
Data Notebook	11
	12

Part Three - The Models

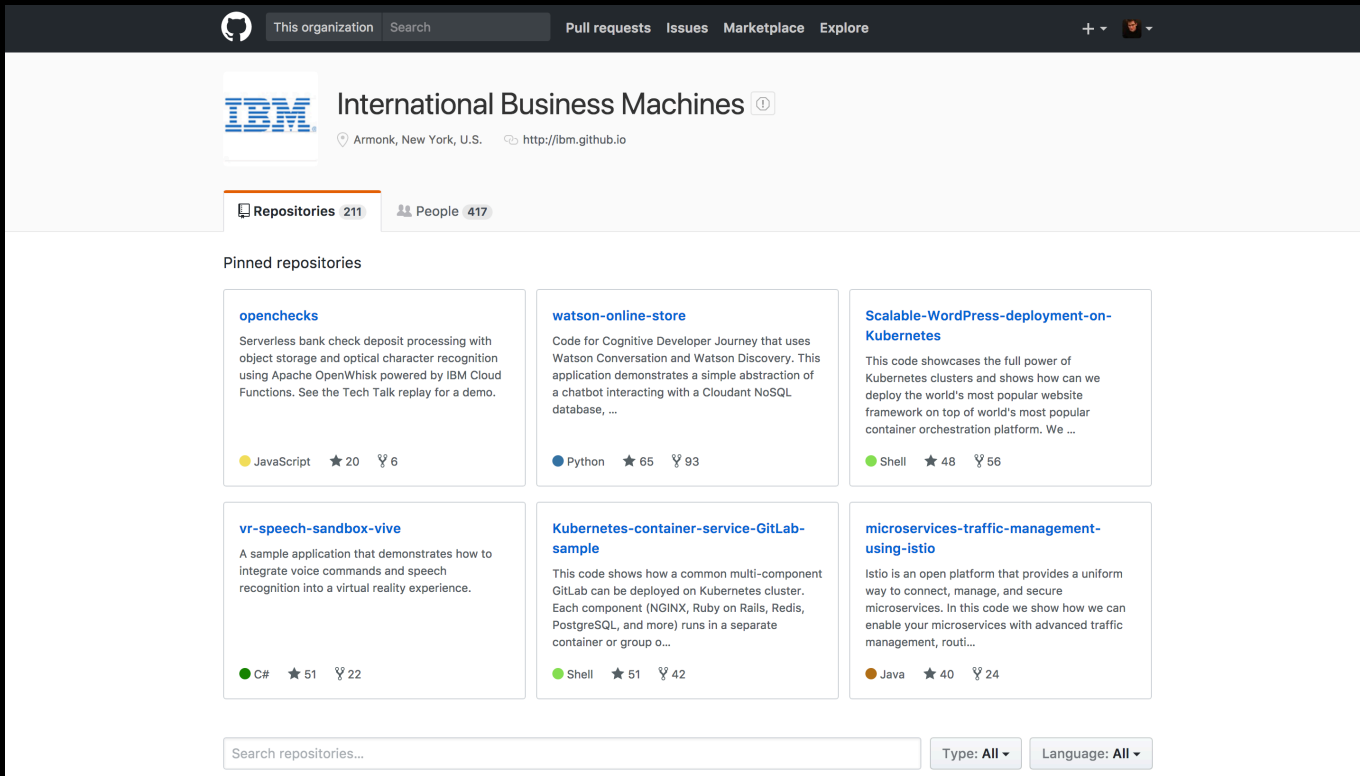
Naïve Bayes Classifier	13
Watson	14
Model Notebook	31
	33

Part Four - The End

Benefits of programming language classification	34
Next Steps	35
Q&A	39
	41

Part 2 The Data

IBM's GitHub



The screenshot displays the GitHub profile for the IBM organization. At the top, the navigation bar includes the GitHub logo, a search bar, and links for Pull requests, Issues, Marketplace, and Explore. The organization's name, 'International Business Machines', is prominently displayed next to its logo, with location and website information below it. A summary bar indicates 211 repositories and 417 people. The 'Pinned repositories' section features six featured projects, each with a brief description, language, star count, and fork count.

Repository Name	Description	Language	Stars	Forks
openchecks	Serverless bank check deposit processing with object storage and optical character recognition using Apache OpenWhisk powered by IBM Cloud Functions. See the Tech Talk replay for a demo.	JavaScript	20	6
watson-online-store	Code for Cognitive Developer Journey that uses Watson Conversation and Watson Discovery. This application demonstrates a simple abstraction of a chatbot interacting with a Cloudant NoSQL database, ...	Python	65	93
Scalable-WordPress-deployment-on-Kubernetes	This code showcases the full power of Kubernetes clusters and shows how can we deploy the world's most popular website framework on top of world's most popular container orchestration platform. We ...	Shell	48	56
vr-speech-sandbox-vive	A sample application that demonstrates how to integrate voice commands and speech recognition into a virtual reality experience.	C#	51	22
Kubernetes-container-service-GitLab-sample	This code shows how a common multi-component GitLab can be deployed on Kubernetes cluster. Each component (NGINX, Ruby on Rails, Redis, PostgreSQL, and more) runs in a separate container or group o...	Shell	51	42
microservices-traffic-management-using-istio	Istio is an open platform that provides a uniform way to connect, manage, and secure microservices. In this code we show how we can enable your microservices with advanced traffic management, routi...	Java	40	24

At the bottom, there is a search bar for repositories and filters for 'Type' (set to All) and 'Language' (set to All).

PyGitHub

```
from github import Github

# First create a Github instance:
g = Github("user", "password")

# Then play with your Github objects:
for repo in g.get_user().get_repos():
    print repo.name
```

plclassifierdata.ipynb

Part 3 The Models

Naïve Bayes Classifier

Bayes Theorem

$$P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)}$$

Bayes Theorem

$$P(\text{programming language} = k | \text{text}) = \frac{P(\text{text} | \text{programming language} = k) P(\text{programming language} = k)}{P(\text{text})}$$

Bayes Theorem

$$P(\text{text} | \text{pl} = k) = ?$$

Naïve Assumption

$$P(\text{word}_1, \text{word}_2, \text{word}_3, \dots | \text{pl} = k) = P(w_1 | \text{pl} = k)^{\#w_1} P(w_2 | \text{pl} = k)^{\#w_2} P(w_3 | \text{pl} = k)^{\#w_3} \dots$$

Naïve Assumption

$$P(\text{text} | \text{pl} = k) = ?$$

Naïve Assumption

$$P(\text{text} | \text{pl} = k) = ?$$

$$P(\text{text} | \text{pl} = k) = P(\text{word}_1, \text{word}_2, \text{word}_3, \dots | \text{pl} = k)$$

Naïve Assumption

$$P(\text{text} | \text{pl} = k) = ?$$

$$\begin{aligned} P(\text{text} | \text{pl} = k) &= P(\text{word}_1, \text{word}_2, \text{word}_3, \dots | \text{pl} = k) \\ &= P(w_1 | \text{pl} = k)^{\#w_1} P(w_2 | \text{pl} = k)^{\#w_2} P(w_3 | \text{pl} = k)^{\#w_3} \dots \end{aligned}$$

Naïve Assumption

$$P(\text{text} | \text{pl} = k) = ?$$

$$\begin{aligned} P(\text{text} | \text{pl} = k) &= P(\text{word}_1, \text{word}_2, \text{word}_3, \dots | \text{pl} = k) \\ &= P(w_1 | \text{pl} = k)^{\#w_1} P(w_2 | \text{pl} = k)^{\#w_2} P(w_3 | \text{pl} = k)^{\#w_3} \dots \end{aligned}$$

$$= \prod_i P(w_i | \text{pl} = k)^{\#w_i}$$

Naïve Bayes Classifier Prediction

$$P(\text{pl} = k | \text{text}) = \frac{P(\text{text} | \text{pl} = k)P(\text{pl} = k)}{P(\text{text})}$$

Naïve Bayes Classifier Prediction

$$P(\text{pl} = k | \text{text}) = \frac{P(\text{text} | \text{pl} = k)P(\text{pl} = k)}{P(\text{text})}$$

Can ignore $P(\text{text})$ as it does not change

$$P(\text{pl} = k | \text{text}) = P(\text{text} | \text{pl} = k)P(\text{pl} = k)$$

Naïve Bayes Classifier Prediction

$$P(\text{pl} = k | \text{text}) = \frac{P(\text{text} | \text{pl} = k) P(\text{pl} = k)}{P(\text{text})}$$

Can ignore $P(\text{text})$ as it does not change

$$\begin{aligned} P(\text{pl} = k | \text{text}) &= P(\text{text} | \text{pl} = k) P(\text{pl} = k) \\ &= \prod_i P(w_i | \text{pl} = k)^{\#w_i} P(\text{pl} = k) \end{aligned}$$

Naïve Bayes Classifier Prediction

$$P(\text{pl} = k | \text{text}) = \prod_i P(w_i | \text{pl} = k)^{\#w_i} P(\text{pl} = k)$$

Naïve Bayes Classifier Prediction

$$P(\text{pl} = k | \text{text}) = \prod_i P(w_i | \text{pl} = k)^{\#w_i} P(\text{pl} = k)$$

$$\log P(\text{pl} = k | \text{text}) = \log \prod_i P(w_i | \text{pl} = k)^{\#w_i} P(\text{pl} = k)$$

Naïve Bayes Classifier Prediction

$$P(\text{pl} = k | \text{text}) = \prod_i P(w_i | \text{pl} = k)^{\#w_i} P(\text{pl} = k)$$

$$\log P(\text{pl} = k | \text{text}) = \log \prod_i P(w_i | \text{pl} = k)^{\#w_i} P(\text{pl} = k)$$

$$= \sum_i \log \#w_i P(w_i | \text{pl} = k) + \log P(\text{pl} = k)$$

Naïve Bayes Classifier Prediction

$$P(\text{pl} = k | \text{text}) = \prod_i P(w_i | \text{pl} = k)^{\#w_i} P(\text{pl} = k)$$

$$\begin{aligned} \log P(\text{pl} = k | \text{text}) &= \log \prod_i P(w_i | \text{pl} = k)^{\#w_i} P(\text{pl} = k) \\ &= \sum_i \log \#w_i P(w_i | \text{pl} = k) + \log P(\text{pl} = k) \end{aligned}$$

$$Y^* = \operatorname{argmax}_k \sum_i \log \#w_i P(w_i | \text{pl} = k) + \log P(\text{pl} = k)$$

Naïve Bayes Classifier Training

$$P(\text{pl} = k | \text{text}) = \prod_i P(w_i | \text{pl} = k)^{\#w_i} P(\text{pl} = k)$$

Naïve Bayes Classifier Training

$$P(\text{pl} = k | \text{text}) = \prod_i P(w_i | \text{pl} = k)^{\#w_i} P(\text{pl} = k)$$

For each programming language:

Probability it occurs in data set (1)

Probability of every word given its in programming language k (2)

$$^1 \frac{\text{count of programs of lang } k}{\text{count of programs}}$$

$$^2 \frac{\text{occurrences of word } i \text{ in pl } k}{\text{words in pl } k}$$

Watson Natural Language Classifier

Watson Natural Language Classifier

Training

- 1 - Call an API

Testing

- 1 - Call an API

plclassifiermodels.ipynb

Part 4 The End

New Developments in Programming Language Classification

GitHub's own Programming Language Classifier, Linguist

Linguist

This library is used on GitHub.com to detect blob languages, ignore binary or vendored files, suppress generated files in diffs, and generate language breakdown graphs.

See [Troubleshooting](#) and [CONTRIBUTING.md](#) before filing an issue or creating a pull request.

Troubleshooting


My repository is detected as the wrong language



The Language stats bar displays languages percentages for the files in the repository. The percentages are calculated based on the bytes of code for each language as reported by the [List Languages](#) API. If the bar is reporting a language that you don't expect:


1. Click on the name of the language in the stats bar to see a list of the files that are identified as that language.
2. If you see files that you didn't write, consider moving the files into one of the [paths for vendored code](#), or use the [manual overrides](#) feature to ignore them.
3. If the files are being misclassified, search for [open issues](#) to see if anyone else has already reported the issue. Any information you can add, especially links to public repositories, is helpful.
4. If there are no reported issues of this misclassification, [open an issue](#) and include a link to the repository or a sample of the code that is being misclassified.

Microsoft's Text Classification Kaggle Competition for Malware Classification



CompetitionsDatasetsKernelsDiscussionJobs...

[Sign In](#)

**Microsoft**

Microsoft Malware Classification Challenge (BIG 2015)

Classify malware into families based on file content and characteristics


\$16,000 · 377 teams · 3 years ago

[Overview](#) [Data](#) [Discussion](#) [Leaderboard](#) [Rules](#)

Overview

[Description](#)
[Evaluation](#)
[Prizes](#)
[Timeline](#)

In recent years, the malware industry has become a well organized market involving large amounts of money. Well funded, multi-player syndicates invest heavily in technologies and capabilities built to evade traditional protection, requiring anti-malware vendors to develop counter mechanisms for finding and deactivating them. In the meantime, they inflict real financial and emotional pain to users of computer systems.



Moving away from DSL in Program Synthesis Research

Published as a conference paper at ICLR 2017

DEEPCODER: LEARNING TO WRITE PROGRAMS

Matej Balog*

Department of Engineering
University of Cambridge

**Alexander L. Gaunt, Marc Brockschmidt,
Sebastian Nowozin, Daniel Tarlow**
Microsoft Research

ABSTRACT

We develop a first line of attack for solving programming competition-style problems from input-output examples using deep learning. The approach is to train a neural network to predict properties of the program that generated the outputs from the inputs. We use the neural network's predictions to augment search techniques from the programming languages community, including enumerative search and an SMT-based solver. Empirically, we show that our approach leads to an order of magnitude speedup over the strong non-augmented baselines and a Recurrent Neural Network approach, and that we are able to solve problems of difficulty comparable to the simplest problems on programming competition websites.

1 INTRODUCTION

A dream of artificial intelligence is to build systems that can write computer programs. Recently, there has been much interest in program-like neural network models (Graves et al., 2014; Weston et al., 2015; Kurach et al., 2015; Joulin & Mikolov, 2015; Grefenstette et al., 2015; Sukhbaatar et al., 2015; Neelakantan et al., 2016; Kaiser & Sutskever, 2016; Reed & de Freitas, 2016; Zaremba et al., 2016; Graves et al., 2016), but none of these can *write programs*; that is, they do not generate human-readable source code. Only very recently, Riedel et al. (2016); Bunel et al. (2016); Gaunt et al. (2016) explored the use of gradient descent to induce source code from input-output examples via differentiable interpreters, and Ling et al. (2016) explored the generation of source code from unstructured text descriptions. However, Gaunt et al. (2016) showed that differentiable interpreter-based program induction is inferior to discrete search-based techniques used by the programming languages community. We are then left with the question of how to make progress on program induction using machine learning techniques.

Next Steps

Last thing

Attachment upload API

Attachment upload API allows you to upload an attachment that you may later send out to many users, without having to repeatedly upload the same data each time it is sent :

```
var attachment = {
  "type": "image",
  "payload": {
    "url": "https://pbs.twimg.com/profile_images/803642201653858305/IAWIDBPw_400x400.png",
    "is_reusable": true
  }
};

controller.api.attachment_upload.upload(attachment, function (err, attachmentId) {
  if(err) {
    // Error
  } else {
    var image = {
      "attachment": {
        "type": "image",
        "payload": {
          "attachment_id": attachmentId
        }
      }
    };
    bot.reply(message, image);
  }
});
```

Q&A

Thank you

Nick Acosta
Developer Advocate

—

nacosta@us.ibm.com
github.com/PubChimps

