# Open Source Software Support

**A Network Tour of Data Science**

Matthaios Olma
Pavlos Nikolopoulos
Stefanos Skalistis

# Adopting Open Source Software (OSS)

Open source software:

- Usage based on some license    - Project lifespan is varying
- Code is fully available

*"wowarmorytools"*

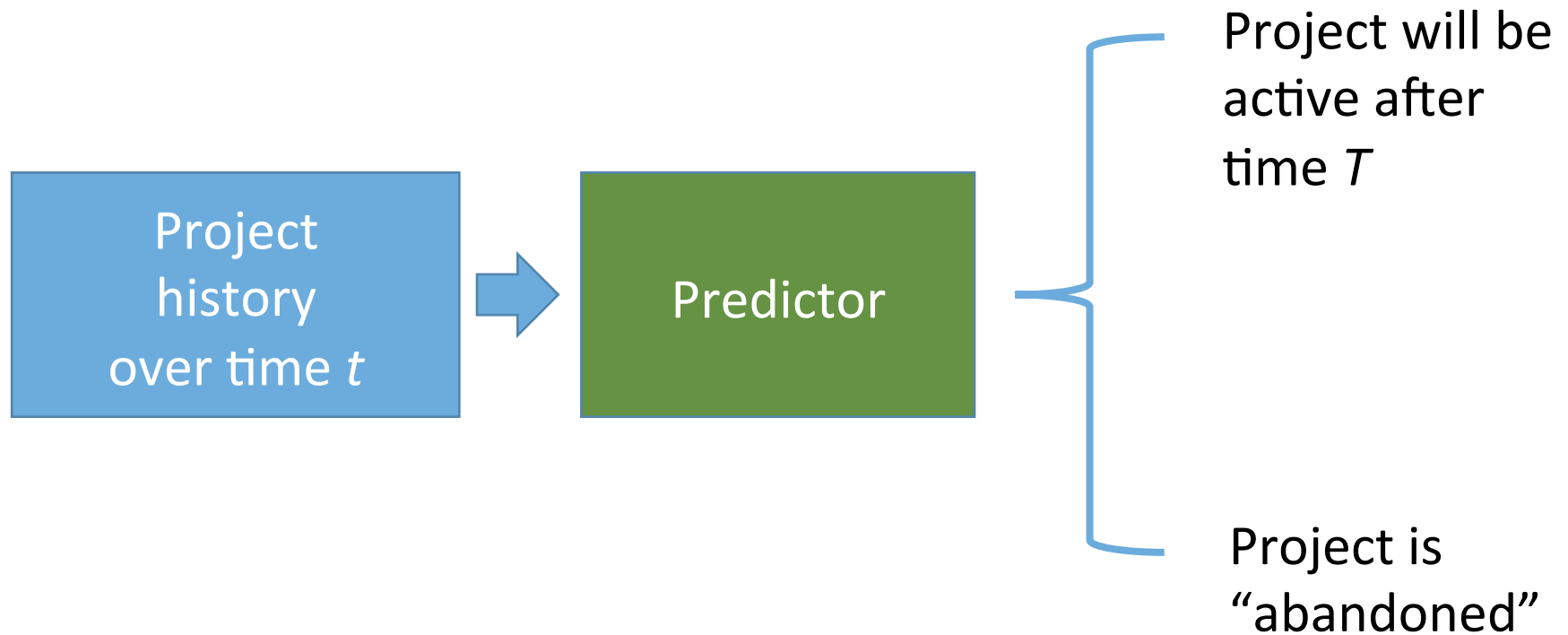*"Blipstick"*

*"James"*

## Should we adopt OSS?

## Will it be supported in the future?

# Predict OSS project survivability

**Desirable Goal:**

*Predict if project will be still active in the future*

Project history over time $t$ → Predictor →

- Project will be active after time $T$
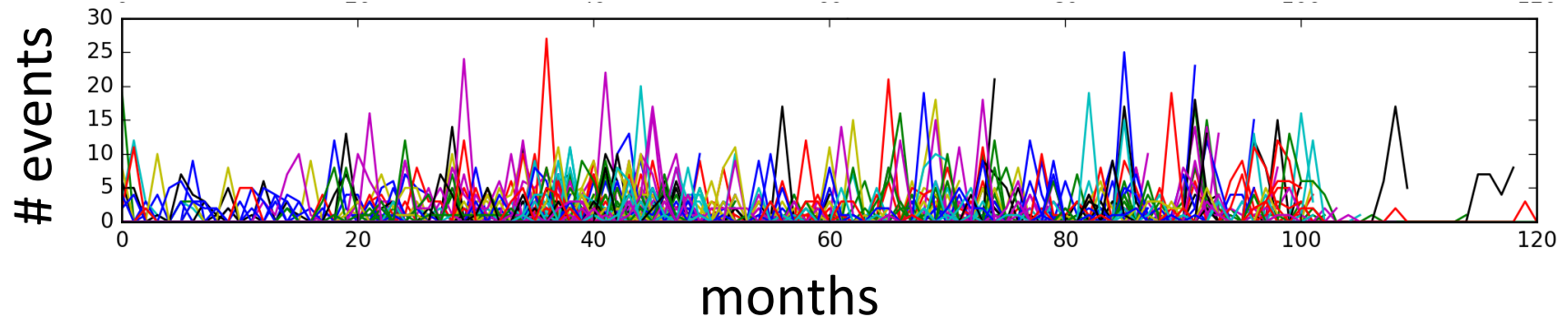- Project is "abandoned"

# Survival prediction is challenging

Success factors:
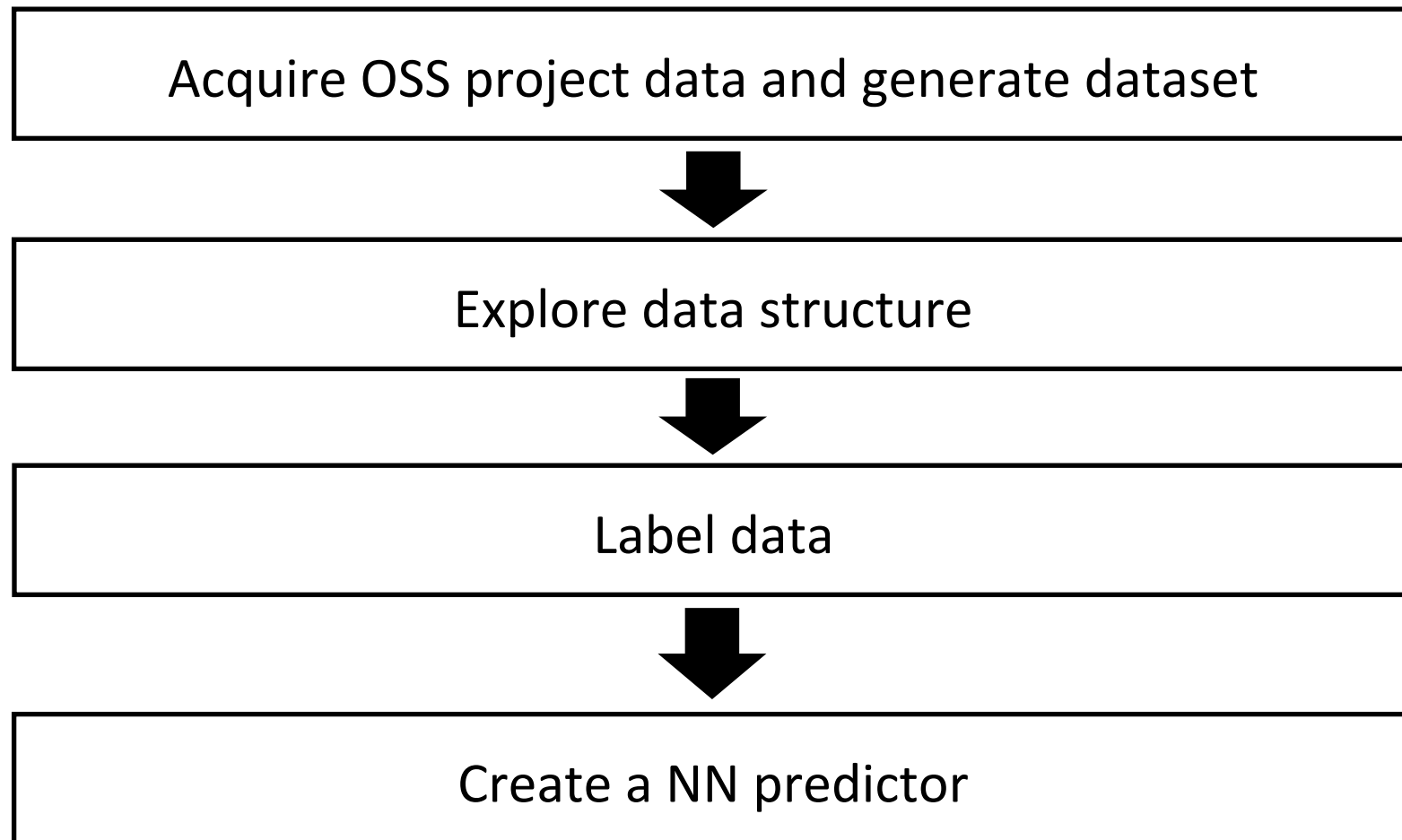- human dynamics
- project popularity
- usefulness

Diverse data:
- Variety of projects
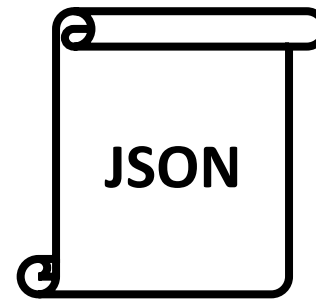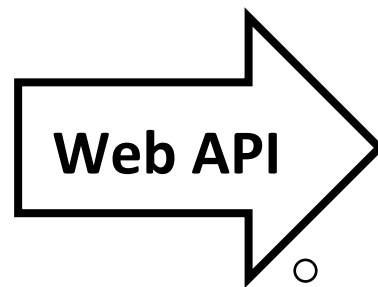- Variety of dev. Techniques
    - *Agile/Waterfall*



**There is no common activity pattern**

# Our Process

Acquire OSS project data and generate dataset

Explore data structure

Label data

Create a NN predictor

# Data Acquisition



3126 projects
Each project is a timeseries of:
- commits
- issues
- comments
- forks
- branches

# Dataset generation

1. Eliminate duplicate projects

2. Create monthly aggregates

<div align="right">

Project #1

$$\begin{pmatrix} 0 & 1 & 0 & 0 & 1 \\ & & \cdots & & \\ 0 & 0 & 1 & 0 & 1 \end{pmatrix}$$

*Nx5*

Project #1

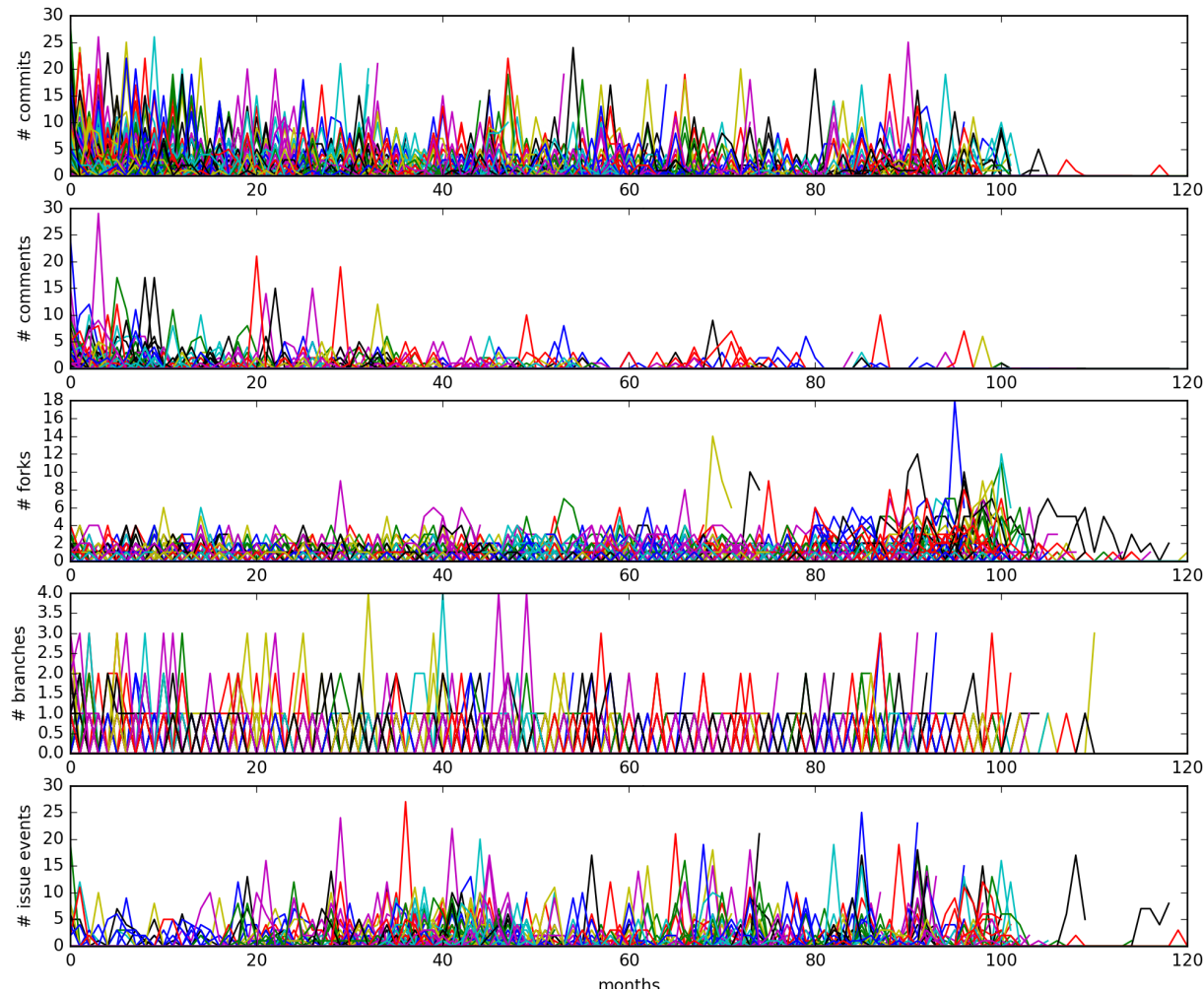$$\begin{pmatrix} 0 & 5 & 0 & 0 & 4 \\ & & \cdots & & \\ 0 & 8 & 0 & 2 & 3 \end{pmatrix}$$

*Mx5*
*M = N/12*

</div>

3. Make all projects start from time 0
   - Start time is the time of the first event

4. Suppress projects with duration < T = 24 months
   - *We are interested in project's activity after time T*

# Data Exploration and Visualization
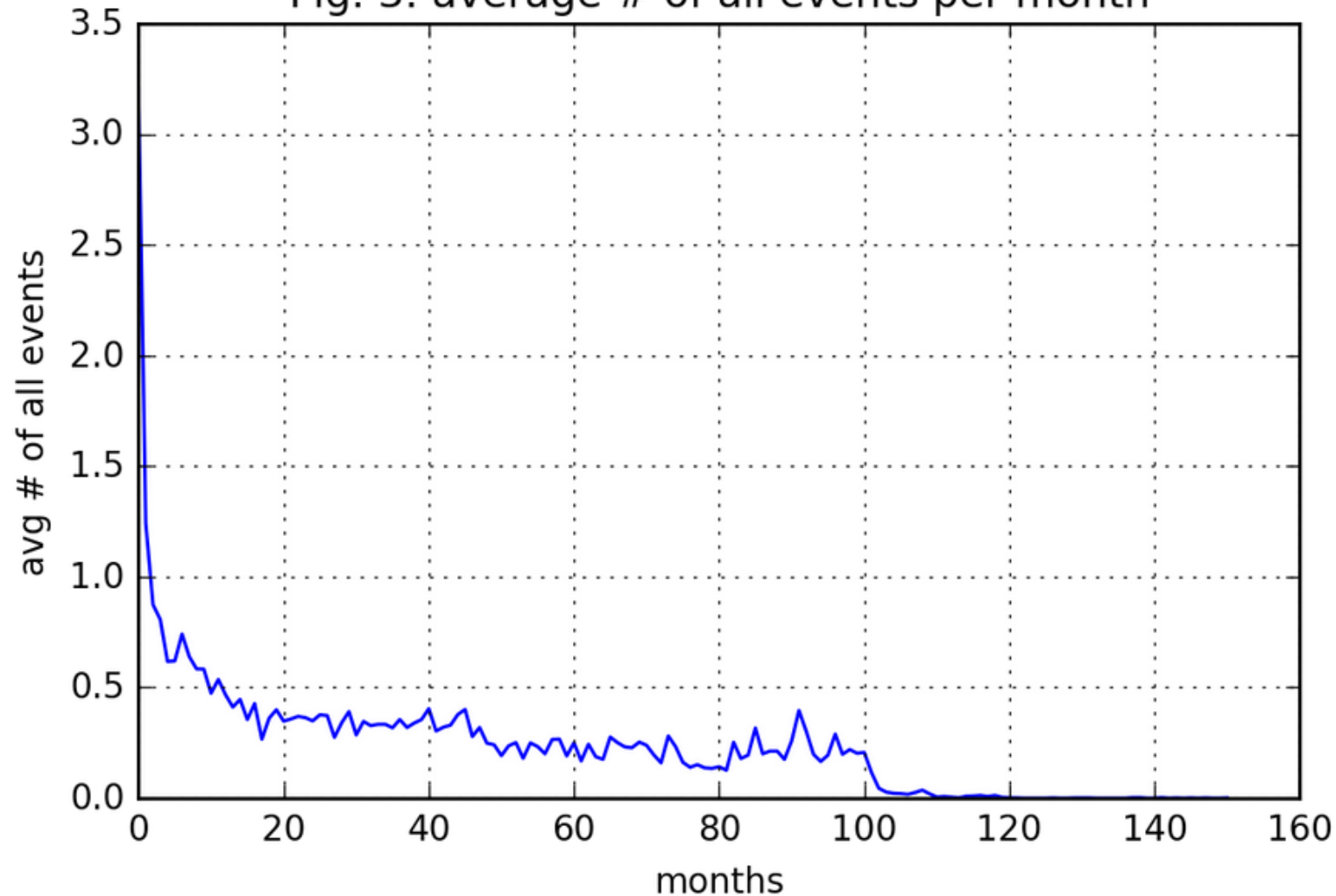
Finally taken into consideration: 834 projects



## no common activity pattern
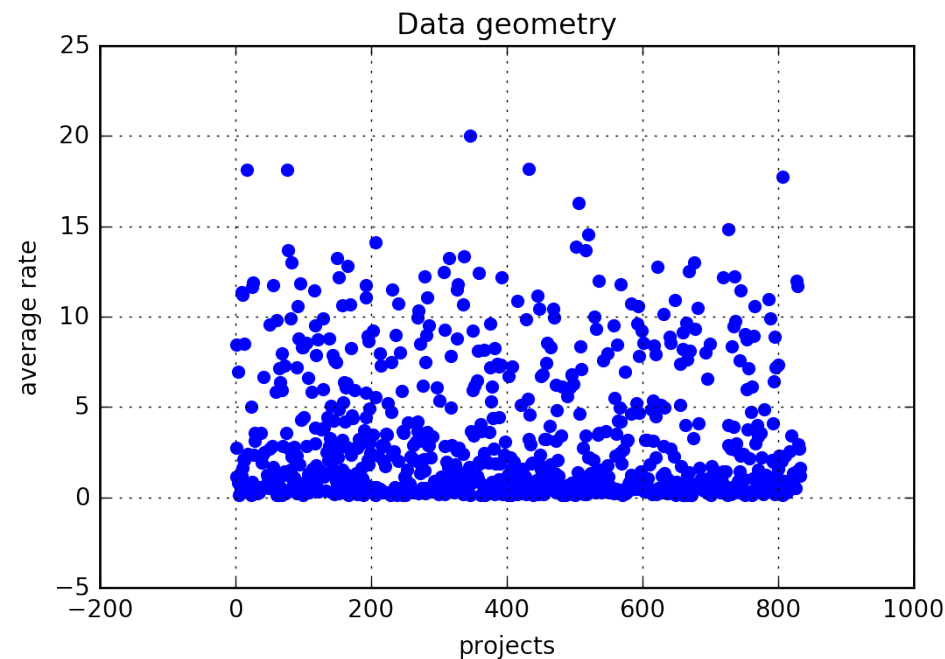
# Choosing the prediction period



Fig. 3: average # of all events per month

**After 24 months the # events converges**

# Data Labeling

- 2 classes: {Active, Inactive}
- Differentiation metric: Average rate of events (after 24 months)
- Threshold set to 6 events/year
  - meaningful threshold in terms of software usage
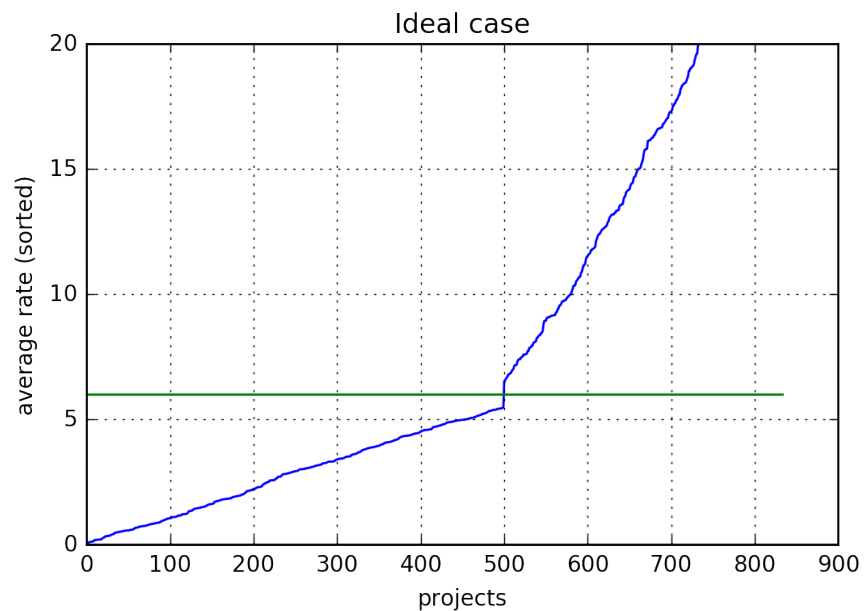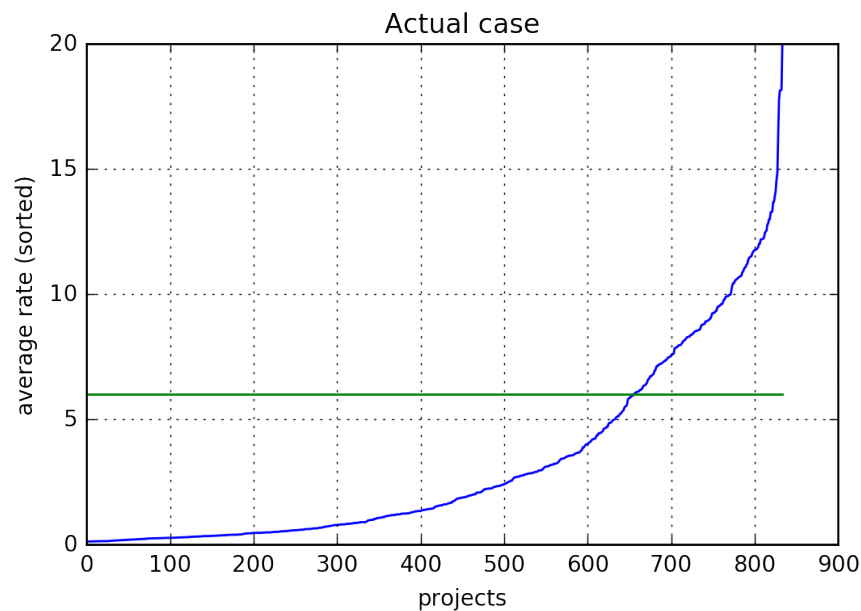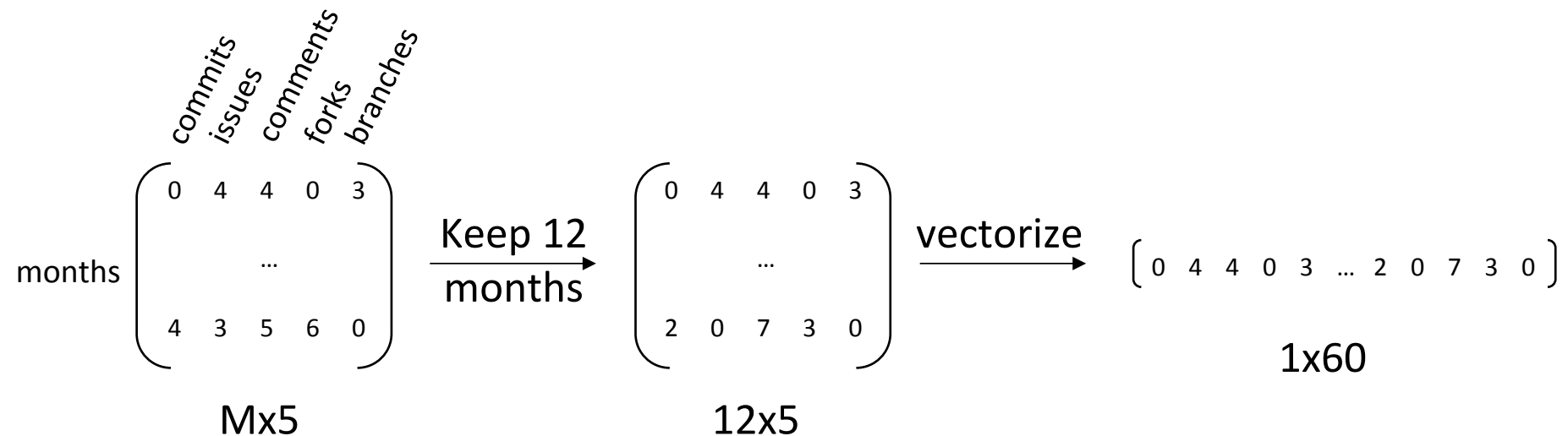  - data is separated in this way



Data geometry

# Data Labeling

- 2 classes: {Active, Inactive}
- Differentiation metric: Average rate of events (after 24 months)
- Threshold set to 6 events/year
  - meaningful threshold in terms of software usage
  - data is separated in this way

# Creating Train and Test sets



- Training set:
  - metadata from the first 12 months of each project
  - 774 projects randomly chosen at batches of 50

- Test set: 50 projects

# 1-layer vs 2-layer NN

**Simple softmax classifier**

$$y = softmax(xW+b)$$

- Normal initialization (std=0.1)
- No regularization

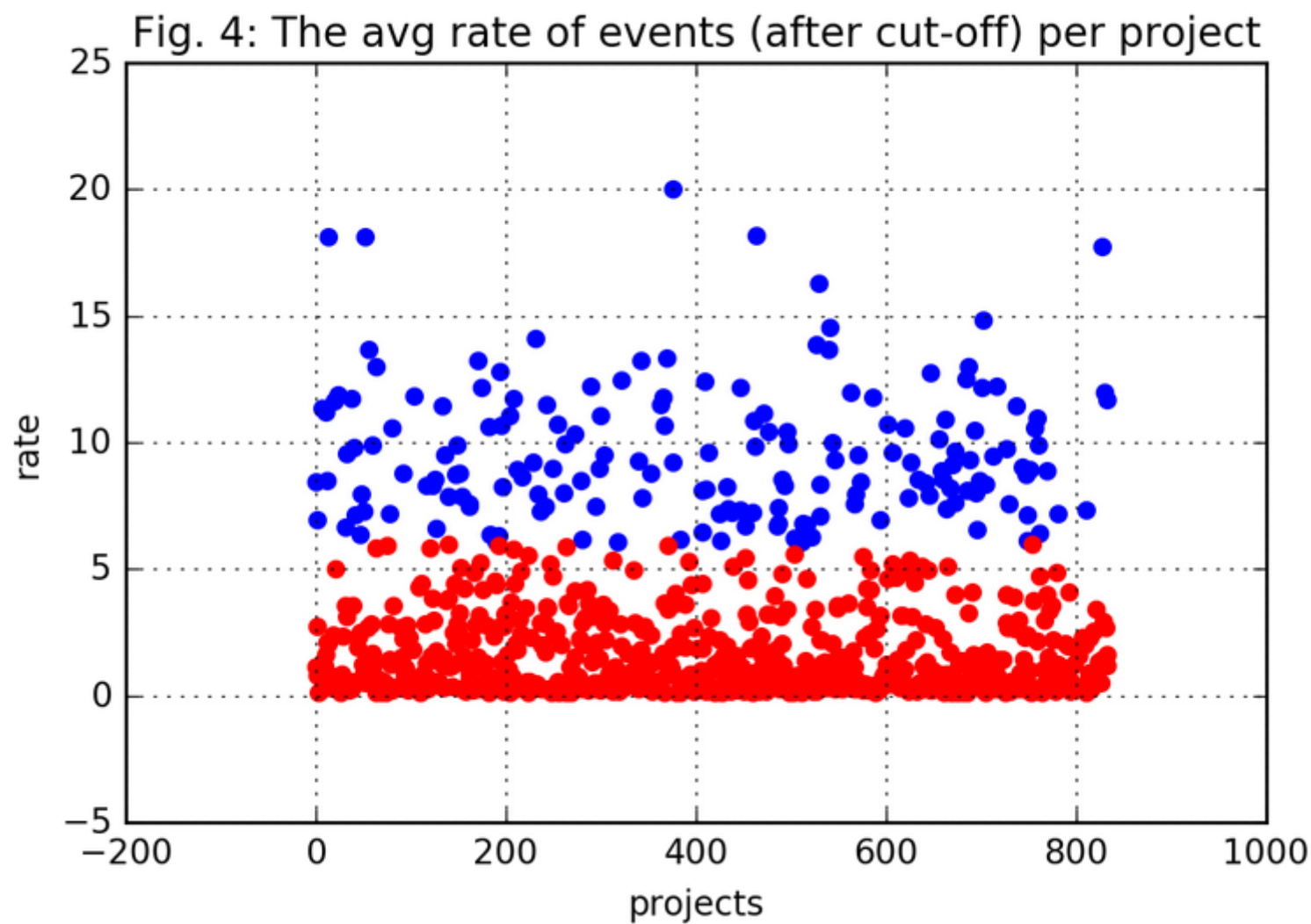- Training accuracy: 76%
- Test accuracy: 80%

**2-layer NN**

$$y = softmax(ReLU(xW_1+b_1)W_2+b_2)$$

- 100 neurons
- Xavier initialization
- ReLU activation
- L2 regularization

- Training accuracy: 92%
- Test accuracy: 84%

# Conclusions

- ## It is possible to predict a project's activity with high accuracy

  - 12 months of metadata are sufficient

- ## Marginal improvement by 2-layer NN over simple linear classifier

  - Labeling based on a linear separation of data
  - Training data and label are based on #events

# Additional no1



Fig. 4: The avg rate of events (after cut-off) per project

# Additional no2



Fig. 1: Cummulative histogram (CDF-like) of project length