# Predicting Software Development Task Times

Rahul Sane

Springboard Intro to Data Science

# The Data
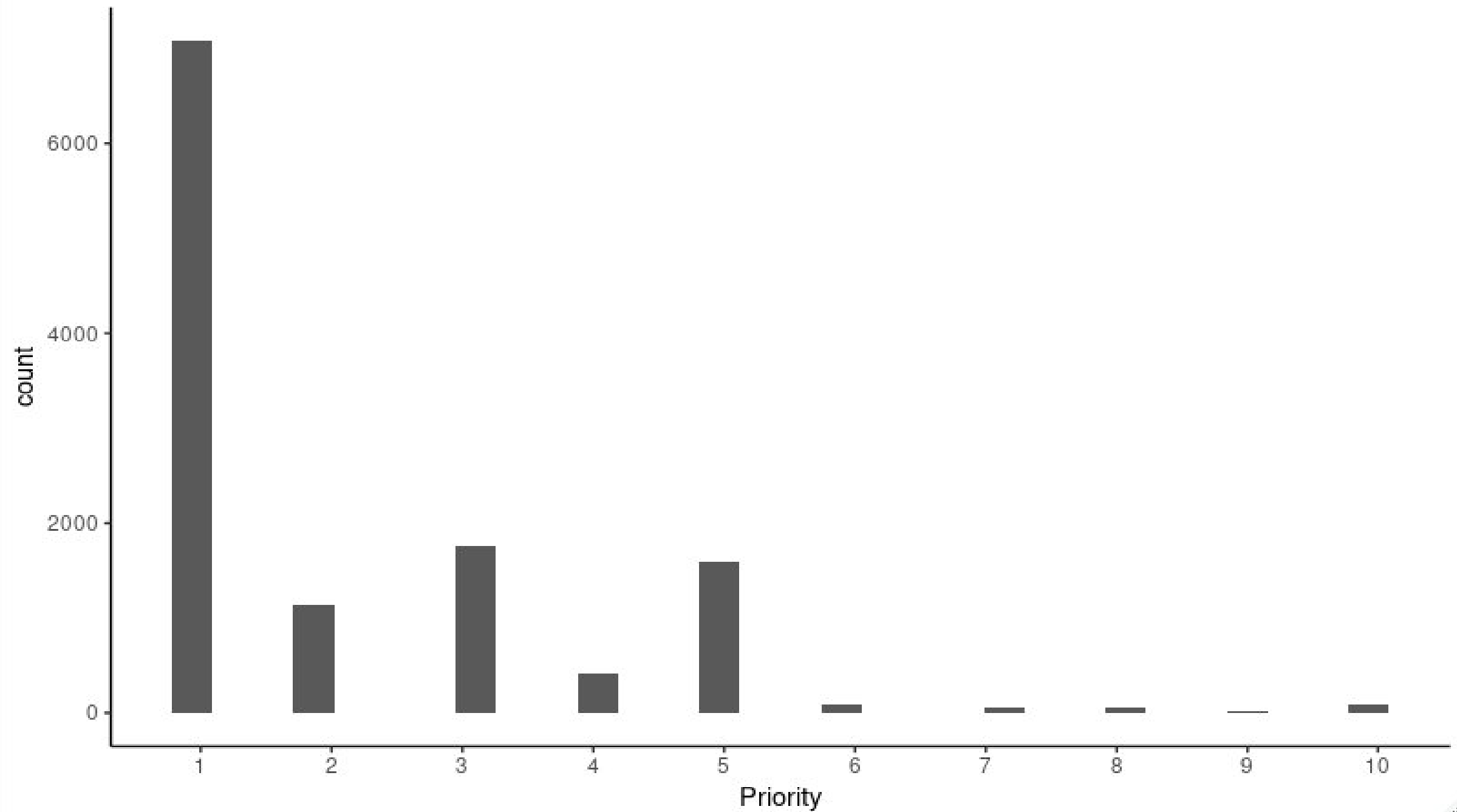
- Historical performance of software developers by task
- Spans 10 years and over 10,000 tasks
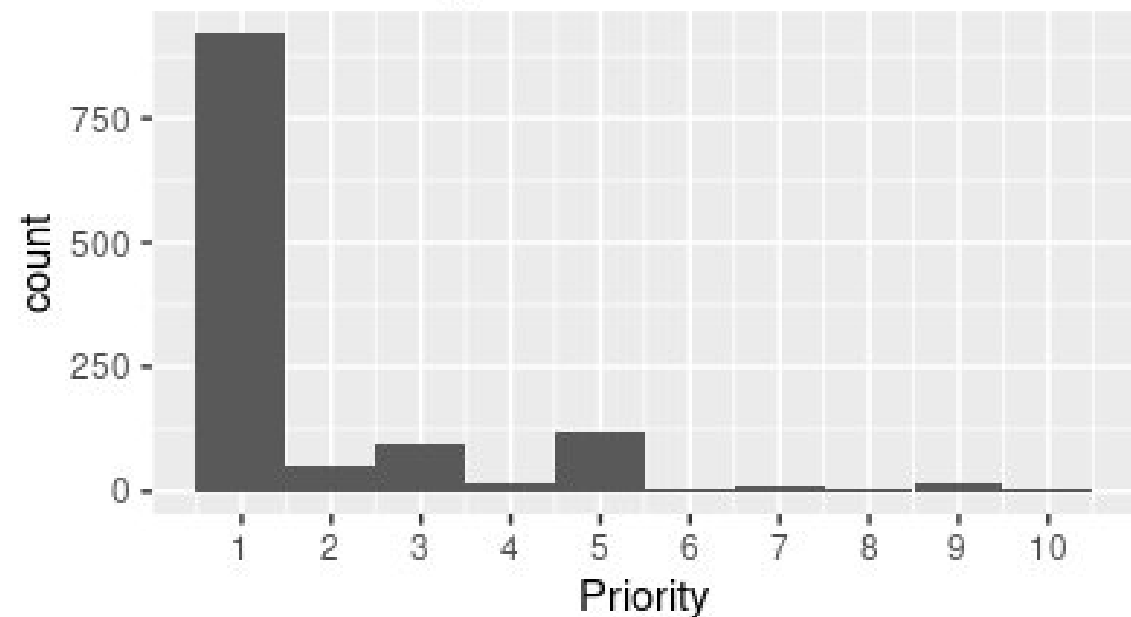- Derek Jones: https://github.com/Derek-Jones/SiP_dataset
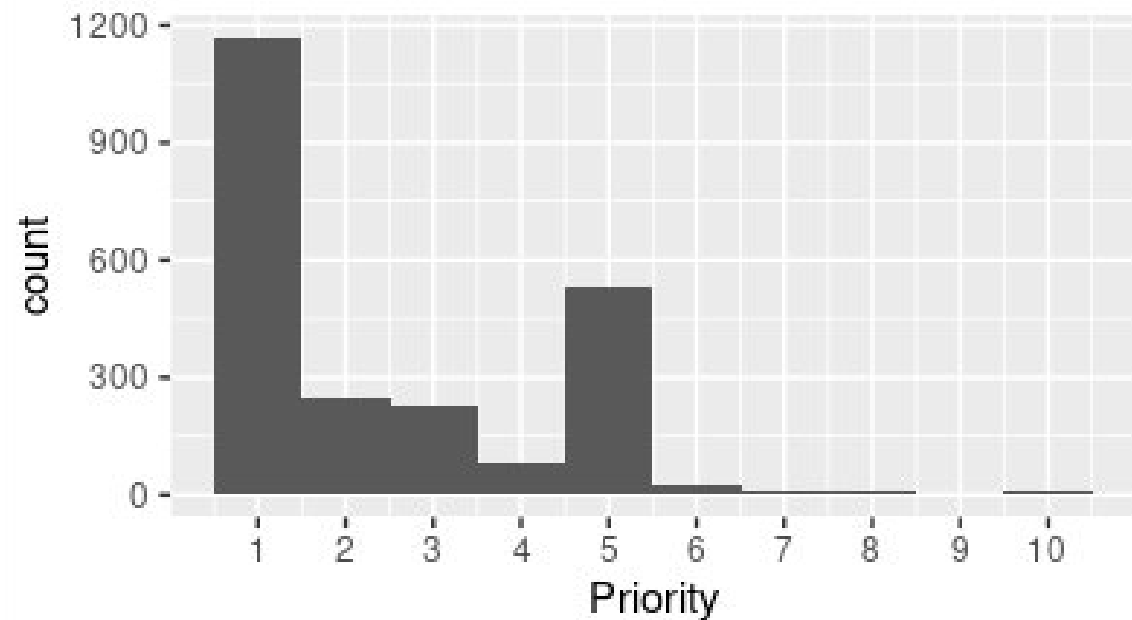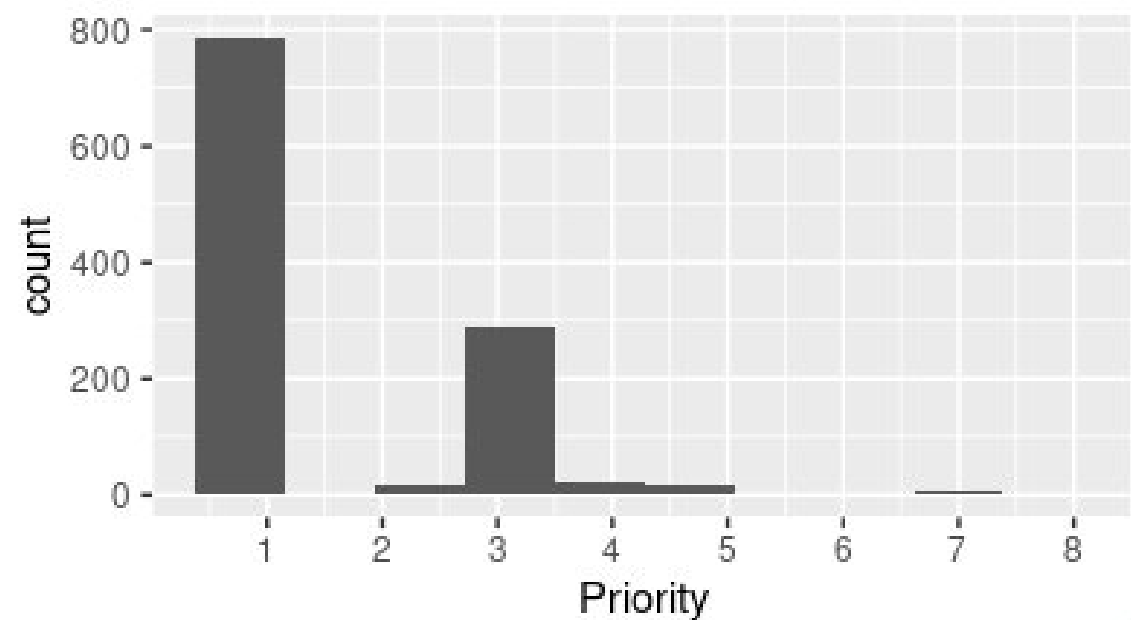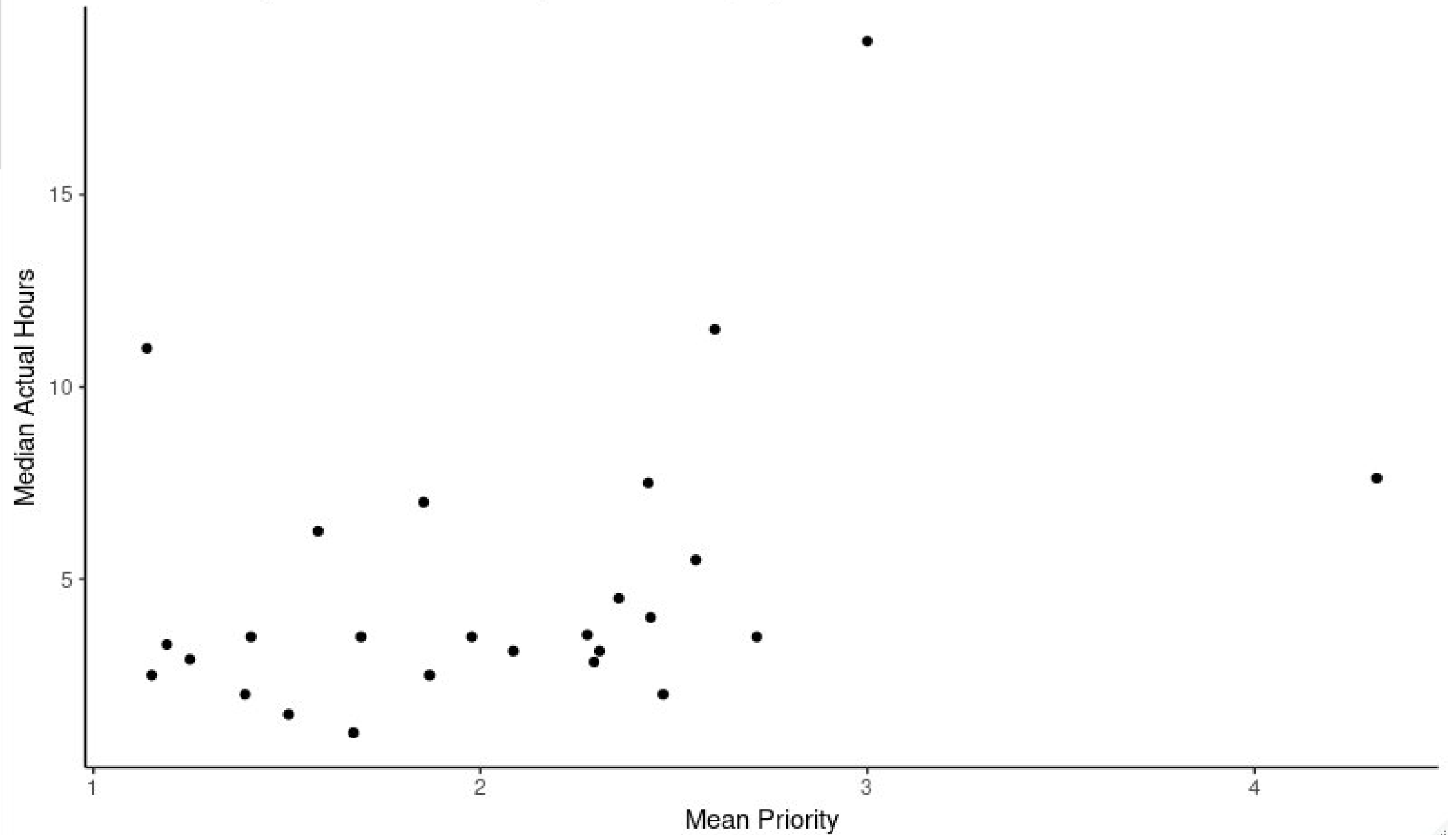
# The Data

- Key columns:
  - TaskNumber
  - Summary
  - Category
  - SubCategory
  - HoursEstimate
  - HoursActual
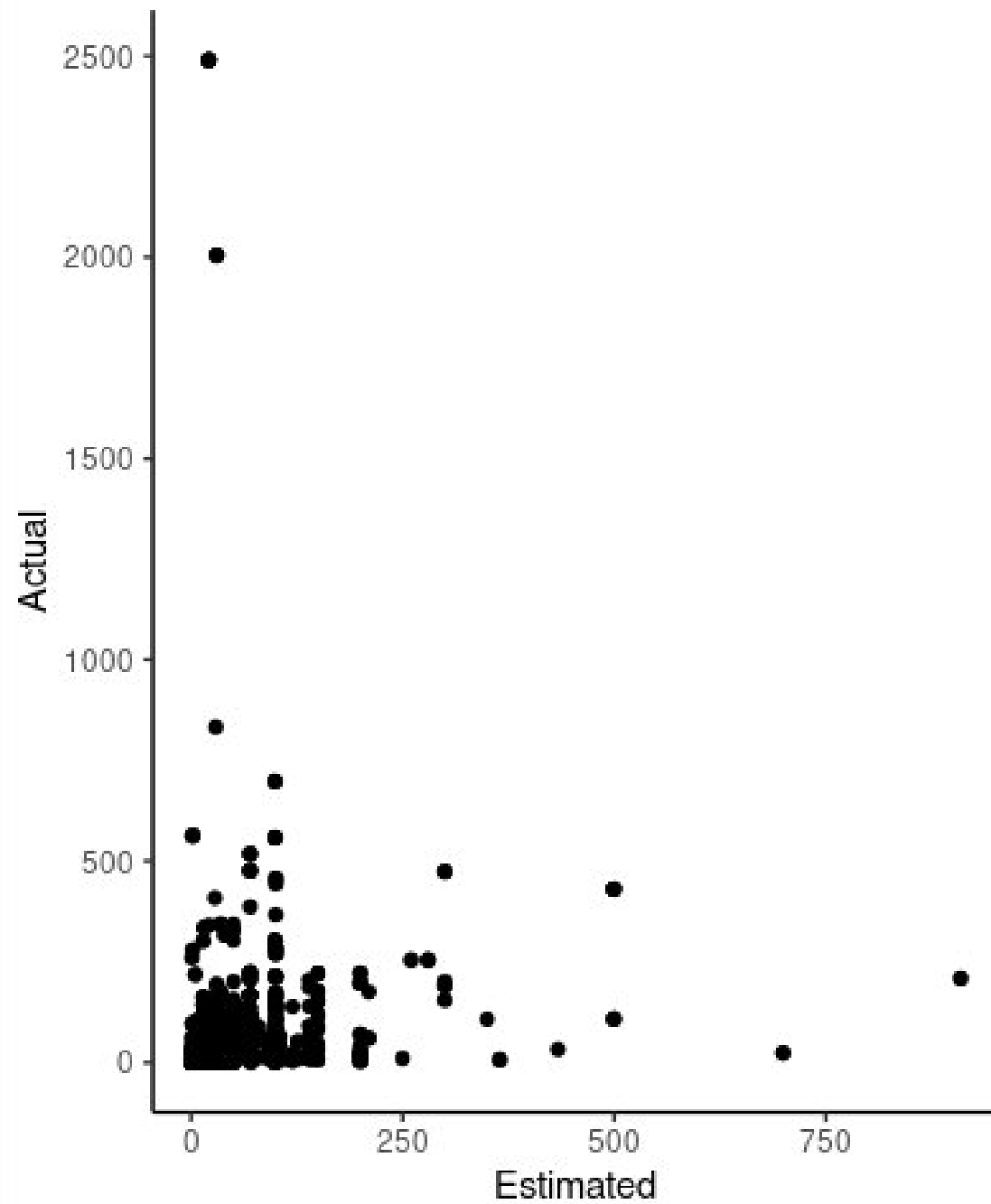  - DeveloperHoursActual

# Distribution of Task Priorities

**A** Enhancement Priorities

**B** In House Support Priorities

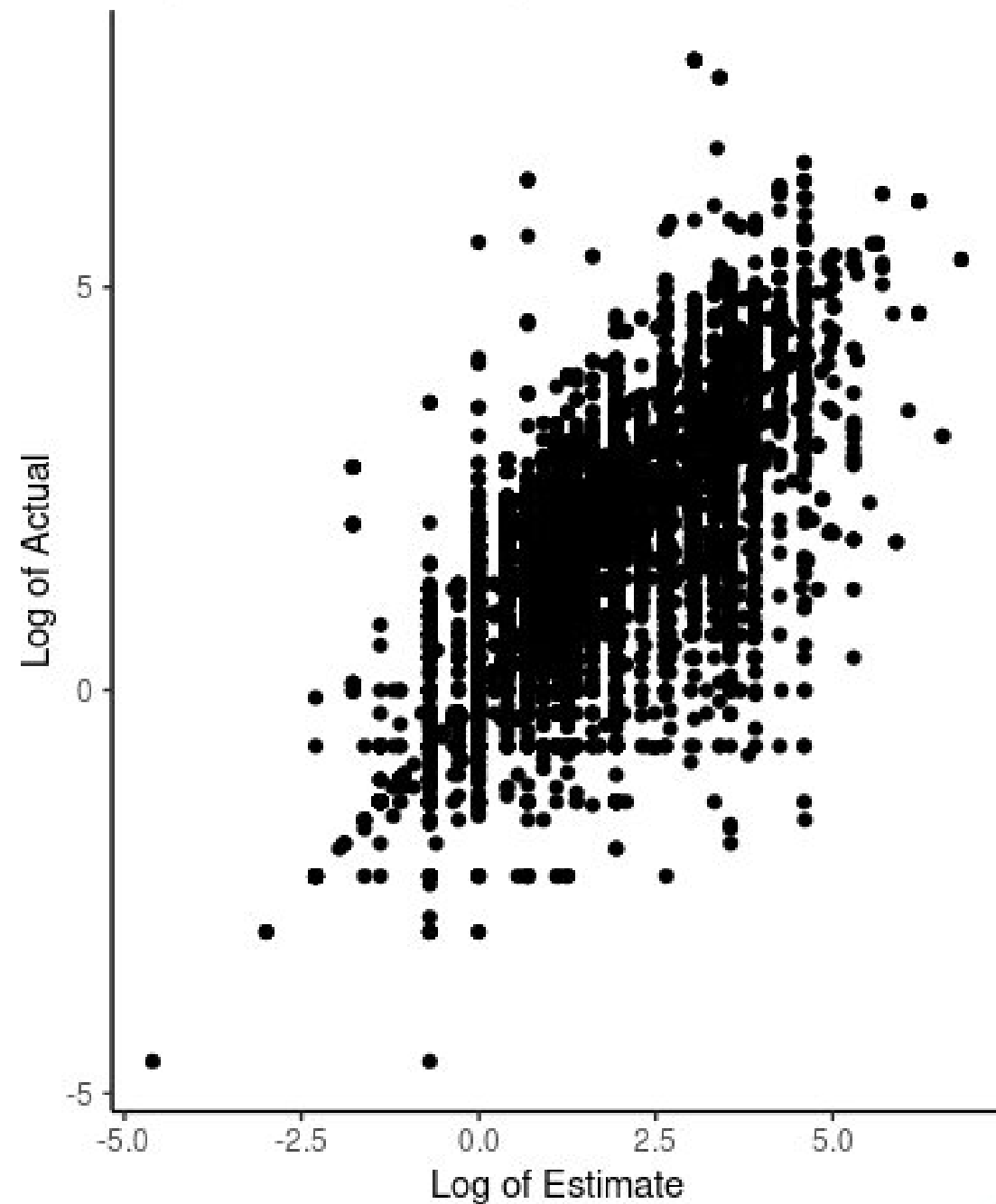**C** Bug Priorities

**D** Support Priorities

Mean Priority vs Median Hours per SubCategory

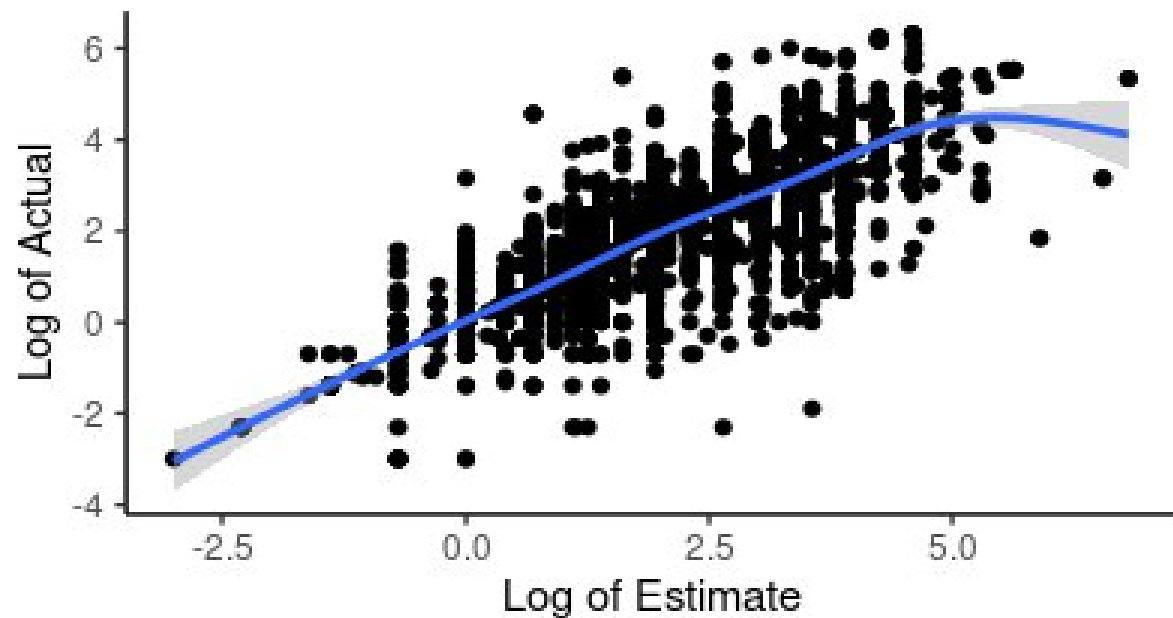**A** Estimated vs Actual Hours

**B** Log-Estimate vs Log-Actual

**A** Enhancement Hours

**B** In House Support Hours

**C** Bug Hours

**D** Support Hours

# Modeling

- Response variable: HoursActual
- Base model using difference between HoursEstimate and HoursActual
- Four machine learning models:
  1. Linear model, HoursEstimate as only predictor
  2. LInear model, all predictors (inc. SubCategory and terms in Summary)
  3. Random forest model, HoursEstimate as only predictor
  4. Random forest model, all predictors

# Model Performance

- Root-mean-square error (RMSE)
- Normalized root-mean-square error (NRMSE); expressed as percentage

# Model Evaluation

| Model | RMSE |
|---|---|
| Base | 27.76 |
| Model 1 | 19.78 |
| Model 2 | 19.85 |
| Model 3 | 20.42 |
| Model 4 | 19.47 |

# Model Evaluation

| Model | RMSE | NRMSE |
|-------|------|-------|
| Model 1 | 19.78 | 5.39% |
| Model 2 | 19.85 | 5.41% |
| Model 3 | 20.42 | 5.57% |
| Model 4 | 19.47 | 5.31% |

# Model Evaluation

- All machine learning models outperform base model
- Simple linear model (Model 2) performed very well
- Model 4 (most complex) is best-performing
  - Most computationally expensive (2.5+ hour runtime)

# Use Case

- Company that relies on internal softwares
- Improve developer workflows
- Improve project management
- Tighter possible timelines

# Further Investigations

- Improve prediction accuracy
- Topic modeling based on summary terms
- Clustering using columns not used here
- Other ML models, inc. SVM and polynomial regression
- Need to avoid over-fitting with more complex models