# Learning From Large Codebases: Chapter 1 (part 2)

# Prediction algorithm

$$y = \underset{y' \in \Omega_x}{\arg\max} \, Pr(y'|x) = \underset{y' \in \Omega_x}{\arg\max} \, score(y', x) = \underset{y' \in \Omega_x}{\arg\max} \, w^T f(y', x)$$

$x$ − input program

$y$ − predicted labels

# Prediction algorithm

1: **for** pass $\in [1 \ldots (\text{num passes})]$ **do**
2:      **for** each node $v$ with unknown property **do**
3:          $E_v \leftarrow$ all edges adjacent to $v$
4:          $\text{score}_v \leftarrow \text{score}(E_v, (y, z))$
5:          **for** $l' \in \text{candidates}(v, (y, z), E_v)$ **do**
6:             $l \leftarrow y_v$
7:             $y_v \leftarrow l'$
8:             **if** $\text{score}(E_v, (y, z)) > score_v$ **and** $y \in \Omega_x$ **then**
9:                 $\text{score}_v \leftarrow \text{score}(E_v, (y, z))$
10:           **else**
11:              $y_v \leftarrow l$

# Prediction algorithm

$$\text{score}(E, A) = \sum_{(a,b,rel) \in E} \sum_{i=1}^{k} w_i \psi_i(A_a, A_b, rel)$$

$$candidates(v, A, E) =$$

$$= \bigcup_{\langle a,v,rel \rangle \in E} \{l^2 \mid \langle l^1, l^2, r \rangle \in topL_s(A_a, rel)\} \quad \cup$$

$$\bigcup_{\langle v,b,rel \rangle \in E} \{l^1 \mid \langle l^1, l^2, r \rangle \in topR_s(A_b, rel)\}$$

# Additional improvements

- Control number of candidates
- Optimization on pairs of nodes

# Learning parameters

Our goal:

$$\forall j, \forall \boldsymbol{y}' \in \Omega_{x^{(j)}} \quad score(\boldsymbol{y}^{(j)}, x^{(j)}) \geq score(\boldsymbol{y}', x^{(j)}) + \Delta(\boldsymbol{y}^{(j)}, \boldsymbol{y}')$$

For every object $\left(x^{(j)}, y^{(j)}\right)$ in training set.

$\Delta(y^{(j)}, y')$ – margin function

# Loss function

Loss function for one object:

$$\ell(\boldsymbol{w}; x^{(j)}, \boldsymbol{y}^{(j)}) = \max_{\boldsymbol{y}' \in \Omega_{x^{(j)}}} \left( \boldsymbol{w}^T [\boldsymbol{f}(\boldsymbol{y}', x^{(j)}) - \boldsymbol{f}(\boldsymbol{y}^{(j)}, x^{(j)})] + \Delta(\boldsymbol{y}^{(j)}, \boldsymbol{y}') \right)$$

Optimal parameters can be found as:

$$\boldsymbol{w}^* = \arg\min_{\boldsymbol{w} \in \mathcal{W}_\lambda} \sum_{j=1}^{t} \ell(\boldsymbol{w}; x^{(j)}, \boldsymbol{y}^{(j)})$$

# Projected Stochastic Gradient Descent

1: $W_\lambda = \{w | w_i \in [0, 1/\lambda] \text{ for all } i\}$
2: **while** $w$ not converged **do**
3:     $g \leftarrow \nabla_w l(w, x^{(j)}, y^{(j)})$
4:     $w \leftarrow w - \alpha g$
5:     $w \leftarrow \text{Proj}_{W_\lambda}(w)$

$$\text{Proj}_{W_\lambda}(w) : w_i' = \max(0, \min(1/\lambda, w_i))$$

# How to compute gradients?

$$g = f(y_{\text{best}}, x^{(j)}) - f(y^{(j)}, x^{(j)})$$

$$y_{\text{best}} = \arg \max_{y' \in \Omega_{x^{(j)}}} (\text{score}(y', x^{(j)}) + \Delta(y^{(j)}, y'))$$
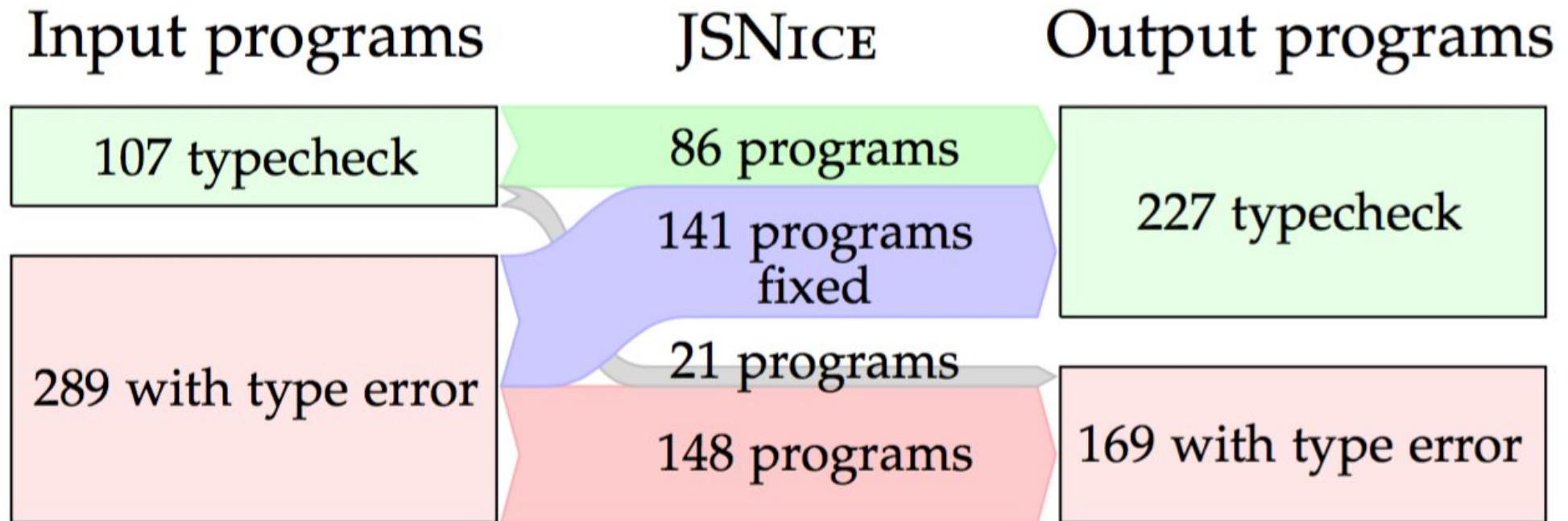
# Details

- Lambda and delta are chosen via cross-validation.
- Train set consists of 324,501 files. Test set consists of 2,710 files.
- Train set is collected from GitHub, test set is collected from BitBucket

# Results

| System | Names Accuracy | Types Precision | Types Recall |
|---|---|---|---|
| **all training data** | **63.4%** | **81.6%** | **66.9%** |
| 10% of training data | 54.5% | 81.4% | 64.8% |
| 1% of training data | 41.2% | 77.9% | 62.8% |
| all data, no structure | 54.1% | 84.0% | 56.0% |
| baseline - no predictions | 25.3% | 37.8% | 100% |

# Typechecking Results

# Running times

| Beam size parameter | Name prediction | | Type prediction | |
| --- | --- | --- | --- | --- |
| $b$ | Accuracy | Time | Precision | Time |
| 4 | 57.9% | 43ms | 80.6% | 36ms |
| 8 | 59.2% | 60ms | 80.9% | 39ms |
| 16 | 62.8% | 62ms | 81.6% | 33ms |
| 32 | 63.2% | 80ms | 81.3% | 37ms |
| 64 **(JSNice)** | 63.4% | 114ms | 81.6% | 40ms |
| 128 | 63.5% | 175ms | 82.0% | 42ms |
| 256 | 63.5% | 275ms | 81.6% | 50ms |
| Naïve greedy, no beam | 62.8% | 115.2 s | 81.7% | 410ms |