# INDEX

Note: Page numbers referring to figures and tables are followed by an italicized *f* or *t* respectively.

persistent malware similarity search systems, *continued*
    commenting on samples, 86
    sample output, 86–87
    searching for similar samples, 86
    wiping database, 86
pick_best_question function, 112–113
pickle module, 143–144
plot function, 162–163, 167
*.png* format, 43
pooling layer, 194
pop instruction, 16–17
Portable Executable (PE) file format, 2–5
    dissecting files using pefile, 5–7
    entry point, 3
        file structure, 2–5, 3*f*
        DOS header, 3
        optional header, 3–4
        PE header, 3
        section headers, 4–5
    sections, defined, 4
Portable Executable (PE) header, 3, 135–136
position independence, 5
precision, 124–126
    effect of base rate on, 124–125
    estimating in deployment environment, 125–126
predict_proba method, 144, 149
PReLU activation function, 179*t*
program stack, defined, 14
projected_graph function, 54
projections, 38
push instruction, 16–17
pyplot module, 148–149, 163

## R

random forest
    overview, 115–116, 116*f*
    random forest–based detectors, 141–146
        complete code for, 144–146
        running detector on new binaries, 144
        training, 142–143
RandomForestClassifier class, 143, 152
ransomware, 30–31, 31f, 155–158, 156*f*, 158, 164–168, 165*f*–166*f*, 168*f*, 172–173, 172*f*–173*f*
.rdata section (in PE file format), 4

Receiver Operating Characteristic curves. *See* ROC curves
rectified linear unit (ReLU) activation function, 177*f*, 178*t*, 180, 182*f*, 183–185, 201
recurrent neural networks (RNNs), 196
registry keys, 32
.reloc section (in PE file format), 5
ReLU (rectified linear unit) activation function, 177*f*, 178*t*, 180, 182*f*, 183–185, 201
ResNets (residual networks), 196–197
resource_projection argument, 52, 227
resource obfuscation, 22
ret instruction, 17–18
reverse engineering, 12
    anti-disassembly techniques, 22
    dynamic analysis for, 26
    methods for, 12
    shared code analysis, 60
    using pefile and capstone, 20–21
RNNs (recurrent neural networks), 196
ROC (Receiver Operating Characteristic) curves, 123–124, 123*f*, 126, 147–150, 230–231, 231*f*
    computing, 147–150
    cross-validation, 151–152, 153*f*
    neural networks, 209–210, 210*f*–211*f*
    visualizing, 149, 150*f*
roc_curve function, 149, 210
.rsrc section (resources) (in PE file format), 4–5

## S

sandbox, 26
Sanders, Hillary, 216
savefig function, 165
scan_file function, 144
scan mode, 230–231
scikit-learn (sklearn) machine learning package, 127–128
    building basic decision tree–based detectors, 129–134
    building random forest–based detectors, 141–146
    evaluating detector performance, 146–153
    feature extraction, 134–135
    hashing trick, 140–141