

JsonGrinder.jl: hierarchical multi-instance data

Explaining HMIL

Šimon Mandlík, Matěj Račinský, Viliam Lisý, *Tomáš Pevný*

Gen digital,
AIC, FEE CTU in Prague

August 21, 2023

Motivation

```
{
  services: [
    {
      port: 22,
      protocol: tcp
    },
    {
      port: 4070,
      protocol: tcp
    },
    {
      port: 4071,
      protocol: tcp
    },
    {
      port: 5353,
      protocol: udp
    }
  ],
  device_id: 8bb8971c-5983-4baa-9753-f0ac21faf162,
  ip: 192.168.1.80,
  mac: ac:63:be:a5:50:43,
  mdns_services: [_workstation._tcp.local., _ssh._tcp.local., _sftp-ssh._tcp.local.]
}
```

Motivation

```
① {  
  device_id: 1fe43d89-329d-40fe-b948-d1cbe0fe6c96,  
  ip: 192.168.15.26,  
  mac: 00:07:25:15:3a:fb,  
  dhcp: [  
    {  
      classid: unknown,  
      paramlist: 1,28,2,3,15,6,119,12,44,47,26,121,42  
    }  
  ]  
}
```

Motivation

```
① {  
  device_id: 1fe43d89-329d-40fe-b948-dicbe0fe6c96,  
  ip: 192.168.15.26,  
  mac: 00:07:25:15:3a:fb,  
  dhcp: [  
    {  
      classid: unknown,  
      paramlist: 1,28,2,3,15,6,119,12,44,47,26,121,42  
    }  
  ]  
}
```

```
② {  
  services: [  
    {  
      port: 5353,  
      protocol: udp  
    },  
    {  
      port: 41800,  
      protocol: tcp  
    }  
  ],  
  device_id: 74ab1b5b-3cb6-4aee-8362-f3b2f016574c,  
  ip: 192.168.0.7,  
  mac: 40:99:22:41:ac:a6,  
  mdns_services: [_spotify-connect._tcp.local.]  
}
```

Motivation

```
① {  
  device_id: 1fe43d89-329d-40fe-b948-dicbe0fe6c96,  
  ip: 192.168.15.26,  
  mac: 00:07:25:15:3a:fb,  
  dhcp: [  
    {  
      classid: unknown,  
      paramlist: 1,28,2,3,15,6,119,12,44,47,26,121,42  
    }  
  ]  
}
```

```
② {  
  services: [  
    {  
      port: 5353,  
      protocol: udp  
    },  
    {  
      port: 41800,  
      protocol: tcp  
    }  
  ],  
  device_id: 74ab1b5b-3cb6-4aee-8362-f3b2f016574c,  
  ip: 192.168.0.7,  
  mac: 40:99:22:41:ac:a6,  
  mdns_services: [_spotify-connect._tcp.local.]  
}
```

```
③ {  
  services: [  
    {  
      port: 80,  
      protocol: tcp  
    },  
    {  
      port: 443,  
      protocol: tcp  
    },  
    {  
      port: 5353,  
      protocol: udp  
    },  
    {  
      port: 5683,  
      protocol: udp  
    }  
  ],  
  device_id: 5347ada9-925c-400e-8a7c-9aedd3c142f6,  
  ip: 192.168.0.147,  
  mac: 80:5e:c0:41:ad:39  
}
```

Motivation

```
① {
  device_id: 1fe43d89-329d-40fe-b948-dicbe0fe6c96,
  ip: 192.168.15.26,
  mac: 00:07:25:15:3a:fb,
  dhcp: [
    {
      classid: unknown,
      paramlist: 1,28,2,3,15,6,119,12,44,47,26,121,42
    }
  ]
}
```

```
② {
  services: [
    {
      port: 5353,
      protocol: udp
    },
    {
      port: 41800,
      protocol: tcp
    }
  ],
  device_id: 74ab1b5b-3cb6-4ae8-8362-f3b2f016574c,
  ip: 192.168.0.7,
  mac: 40:99:22:41:ac:a6,
  mdns_services: [_spotify-connect._tcp.local.]
}
```

```
③ {
  services: [
    {
      port: 80,
      protocol: tcp
    },
    {
      port: 443,
      protocol: tcp
    },
    {
      port: 5353,
      protocol: udp
    },
    {
      port: 5683,
      protocol: udp
    }
  ],
  device_id: 5347ada9-925c-400e-8a7c-9aedd3c142f6,
  ip: 192.168.0.147,
  mac: 80:5e:c0:41:ad:39
}
```

```
④ {
  services: [
    {
      port: 80,
      protocol: tcp
    },
    {
      port: 80,
      protocol: tcp
    },
    {
      port: 111,
      protocol: tcp
    },
    {
      port: 111,
      protocol: tcp
    }
  ],
  device_id: unknown,
  location: http://192.168.1.121:60006/upmp/desc/aio_device/aio_device.xml,
  st: .,
  st: urn:schemas-demon-com:device:AioServices:1,
  server: LINUX UPnP/1.0 Demon-Haos/139046,
  user_agent: unknown
},
{
  location: http://192.168.1.121:60006/upmp/desc/aio_device/aio_device.xml,
  st: .,
  st: uuid:bf87346f-27aa-1691-0080-0005cda2625e,
  server: LINUX UPnP/1.0 Demon-Haos/139046,
  user_agent: unknown
},
{
  location: http://192.168.1.121:60006/upmp/desc/aio_device/aio_device.xml,
  st: .,
  st: uuid:5cdaa5bc-2d82-1022-0080-0005cda2625e,
  server: LINUX UPnP/1.0 Demon-Haos/139046,
  user_agent: unknown
},
{
  location: http://192.168.1.121:60006/upmp/desc/aio_device/aio_device.xml,
  st: .,
  st: urn:schemas-upmp-org:service:AVTransport:1,
  server: LINUX UPnP/1.0 Demon-Haos/139046,
  user_agent: unknown
},
{
  location: http://192.168.1.121:60006/upmp/desc/aio_device/aio_device.xml,
  st: .,
  st: urn:schemas-upmp-org:device:MediaRenderer:1,
  server: LINUX UPnP/1.0 Demon-Haos/139046,
  user_agent: unknown
},
{
  location: http://192.168.1.121:60006/upmp/desc/aio_device/aio_device.xml,
  st: .,
  st: urn:schemas-demon-com:device:MediaServer:1,
  server: LINUX UPnP/1.0 Demon-Haos/139046,
  user_agent: unknown
},
  services: [urn:demon-com:serviceId:ACT]
},
{
  model_name: Demon AVS-X2408,
  manufacturer: Demon,
  device_type: urn:schemas-upmp-org:device:MediaServer:1,
  model_description: Shares User defined folders and files to other Universal Plug and Play media devices.,
  services: [urn:upmp-org:serviceId:ContentDirectory, urn:upmp-org:serviceId:ConnectionManager]
}],
  mdns_services: [_airplay._tcp.local., _spotify-connect._tcp.local., _sop._tcp.local., _http._tcp.local.]
}
```

Our options

- Feature engineering
- JSON as an image
- JSON as text
- Heterogenous graph
- HMIL

HMIL's design

```
{
  services: [
    {
      port: 22,
      protocol: tcp
    },
    {
      port: 2071,
      protocol: tcp
    },
    {
      port: 5353,
      protocol: udp
    }
  ],
  ip: 192.168.1.80,
  mac: ac:63:be:a5:50:43,
  mdns_services: [
    _workstation._tcp.local.,
    _ssh._tcp.local.,
    _sftp-ssh._tcp.local.
  ]
}
```

Project each node to vector
assuming all his childs are vectors.

Leafs — numbers

```
{
  services: [
    {
      port: 22, —————> □
      protocol: tcp
    },
    {
      port: 2071, —————> □
      protocol: tcp
    },
    {
      port: 5353, —————> □
      protocol: udp
    }
  ],
  ip: 192.168.1.80,
  mac: ac:63:be:a5:50:43,
  mdns_services: [
    _workstation._tcp.local.,
    _ssh._tcp.local.,
    _sftp-ssh._tcp.local.
  ]
}
```

Leafs — categorical features

```
{
  services: [
    {
      port: 22,
      protocol: tcp
    },
    {
      port: 2071,
      protocol: tcp
    },
    {
      port: 5353,
      protocol: udp
    }
  ],
  ip: 192.168.1.80,
  mac: ac:63:be:a5:50:43,
  mdns_services: [
    _workstation._tcp.local.,
    _ssh._tcp.local.,
    _sftp-ssh._tcp.local.
  ]
}
```

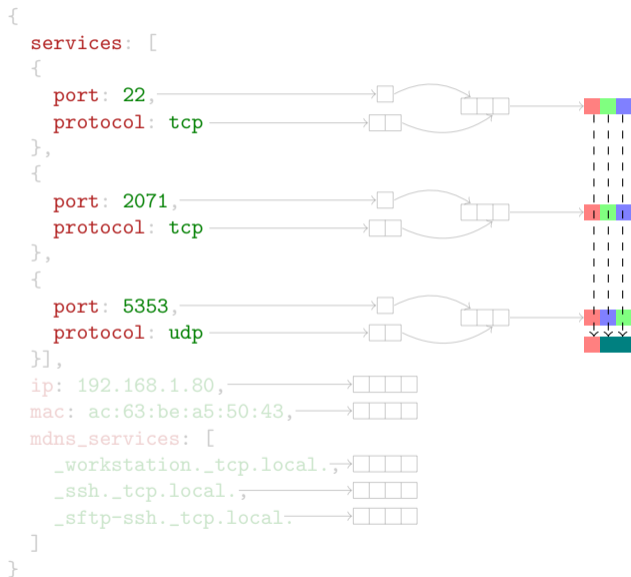
Leafs — string features

```
{
  services: [
    {
      port: 22, —————> □
      protocol: tcp —————> □□
    },
    {
      port: 2071, —————> □
      protocol: tcp —————> □□
    },
    {
      port: 5353, —————> □
      protocol: udp —————> □□
    }
  ],
  ip: 192.168.1.80, —————> □□□□
  mac: ac:63:be:a5:50:43, —————> □□□□
  mdns_services: [
    _workstation._tcp.local., —————> □□□□
    _ssh._tcp.local., —————> □□□□
    _sftp-ssh._tcp.local. —————> □□□□
  ]
}
```

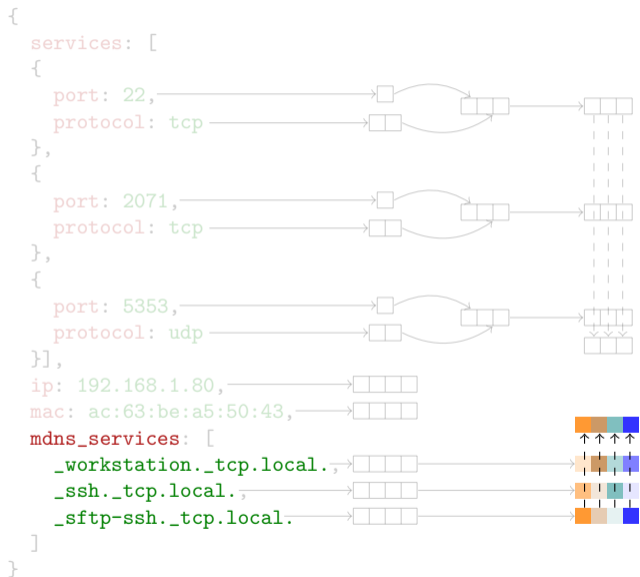
Dictionaries of ports and protocols

```
{  
  services: [  
    {  
      port: 22, —————→ [red]  
      protocol: tcp —————→ [green][blue]  
    },  
    {  
      port: 2071, —————→ [red]  
      protocol: tcp —————→ [green][blue]  
    },  
    {  
      port: 5353, —————→ [red]  
      protocol: udp —————→ [blue][green]  
    }  
  ],  
  ip: 192.168.1.80, —————→ [ ][ ][ ][ ]  
  mac: ac:63:be:a5:50:43, —————→ [ ][ ][ ][ ][ ][ ]  
  mdns_services: [  
    _workstation._tcp.local., —————→ [ ][ ][ ][ ]  
    _ssh._tcp.local., —————→ [ ][ ][ ][ ]  
    _sftp-ssh._tcp.local., —————→ [ ][ ][ ][ ]  
  ]  
}
```

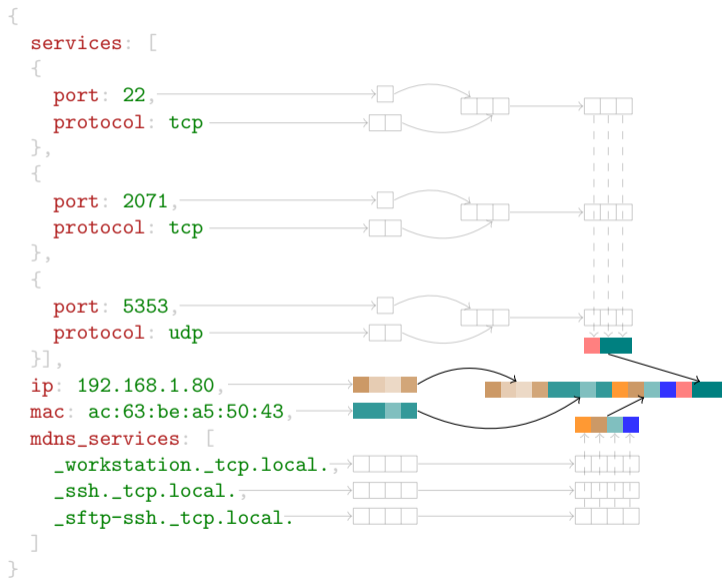
Array of services



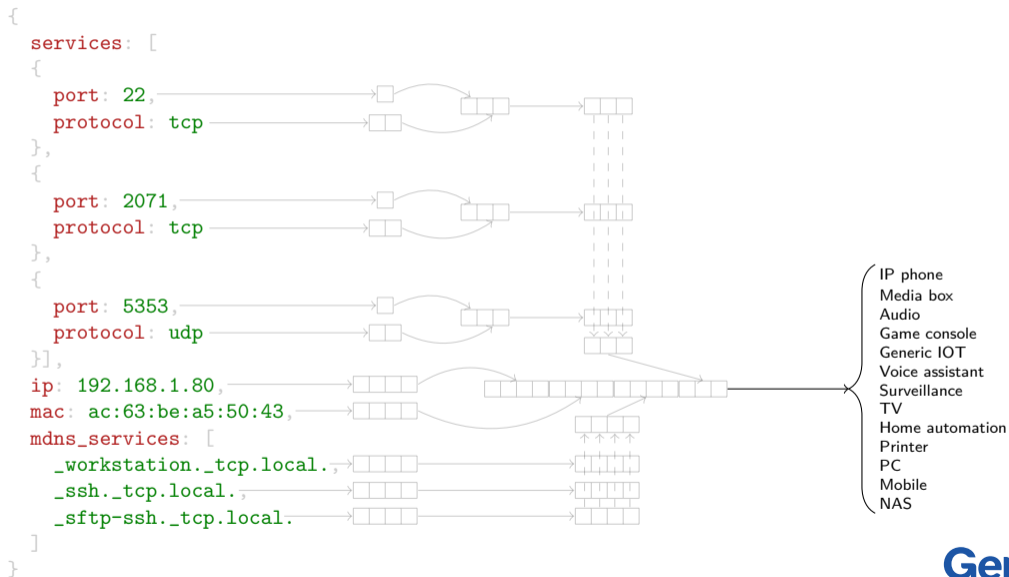
Array of mdns services



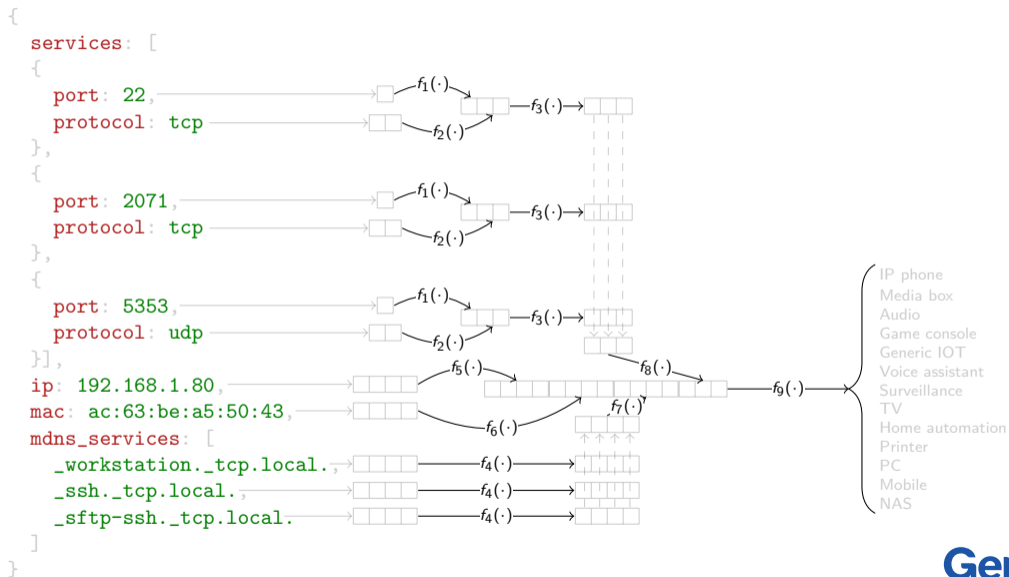
Topmost dictionary



Predicting classes



Complete model



Example of use

```
Code Blame 29 lines (22 loc) · 960 Bytes Raw Copy Download Edit View Source
```

```
1 using JsonGrinder, Mill, Flux, MLDatasets, JSON, CSV, DataFrames, Statistics, Random
2
3 Random.seed!(42)
4
5 include("prepare_data.jl")
6
7 BATCH_SIZE = 50
8
9 x_train, y_train = load_train_data()
10 schema = JsonGrinder.schema(x_train)
11 extractor = suggestextractor(schema)
12 ds_train = map(extractor, x_train)
13
14 encoder = reflectinmodel(schema, extractor, d -> Chain(Dense(d, 50, relu), Dropout(0.25)))
15 model = Dense(50, length(LABELS)) * encoder
16 loss = (ds, y) -> Flux.Losses.logitcrossentropy(model(ds), Flux.onehotbatch(y, eachindex(LABELS)))
17 accuracy = (ds, y) -> mean(Flux.onecold(model(ds), eachindex(LABELS)) .== y)
18
19 opt = AdaBelief()
20 ps = Flux.params(model)
21 data_loader = Flux.Data.DataLoader(ds_train, y_train, batchsize=BATCH_SIZE, shuffle=true)
22 Flux.@epochs 3 begin
23     Flux.Optimise.train!(loss, ps, data_loader, opt)
24     @show accuracy(ds_train, y_train)
25 end
26
27 x_test, y_test = load_test_data()
28 ds_test = map(extractor, x_test)
29 @show accuracy(ds_test, y_test)
```

Comparison to SOTA

Dataset	Sample Size	Accuracy		
		Default	Tuned	Competing
Mutagenesis	4k-8k	0.886	0.932	0.912
Device ID	0.1k-0.3M	0.932	0.971	0.967
EMBER	3k-6M	0.948	0.996	0.974

Try it



<https://github.com/CTUAvastLab/JsonGrinder.jl>