

# Unified Streaming/Batch Learning and Explainable Multi-output Prediction

Jesse Read



- 1 Data Streams as Time Series
- 2 Chaining Methods for Multi-Output Learning
- 3 Applications of Chaining in Data Streams

# Outline

- 1 Data Streams as Time Series
- 2 Chaining Methods for Multi-Output Learning
- 3 Applications of Chaining in Data Streams

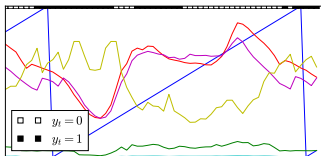
# Data Streams

A data stream,

$$x_1, x_2, \dots, x_t, \dots$$

where, at real time  $t$  we observe  $x_t$ , which comes from some **concept** (which we don't observe directly):

$$x_t \sim P_t$$



Date	Time	Dir	Remote IP Addr	Remote Name / Message	R Port	Local IP Addr	L Port
07/26	15:37:54.01	o	udp 48.13.99.178	48-10-99-178.bto-net.bg	26381	192.168.1.117	29011
07/26	15:37:54.01	o	udp 190.22.141.163	190-22-141-163.baf.movistar.cl	28431	192.168.1.117	29011
07/26	15:37:54.01	o	udp 193.160.237.215		137	192.168.1.117	137
07/26	15:37:54.01	o	udp 55.136.9.212	asbl-dynamic-55-136-9-212.galcsinfo.net	30837	192.168.1.117	29011
07/26	15:37:49.34	o	udp 85.138.197.48	asbl-138-197-48.spe.netcabo.pt	42093	192.168.1.117	29011
07/26	15:37:49.34	o	udp 195.190.109.190	asbl-195-190-109-190.sovintnet.ru	19705	192.168.1.117	29011
07/26	15:37:49.34	o	udp 173.76.100.81	post-173-76-100-81.tampfl.fox.verizon.net	67140	192.168.1.117	29011
07/26	15:37:49.34	o	udp 92.37.48.91	isp-92-37-48-91.dynamic.asia.net	30537	192.168.1.117	29011
07/26	15:37:49.34	o	udp 109.70.188.205		137	192.168.1.117	137
07/26	15:37:43.77	o	udp 178.73.102.9		137	192.168.1.117	137
07/26	15:37:43.77	o	udp 190.229.28.92	host92.190-229-28.telecom.net.br	137	192.168.1.117	137
07/26	15:37:43.77	o	udp 89.73.245.180	89-73-245-180.dynamic.chello.it	137	192.168.1.117	137
07/26	15:37:43.77	o	udp 89.28.31.190		137	192.168.1.117	137
07/26	15:37:43.77	o	udp 10.178.64.1		68	192.168.1.117	68
07/26	15:37:43.77	o	udp 172.18.41.9		255	192.168.1.117	255
07/26	15:37:43.77	o	udp 84.233.291.170		117	192.168.1.117	29011
07/26	15:37:38.00	o	udp 79.113.211.58	79-113-211-58.ednet.ro	1024	192.168.1.117	29011



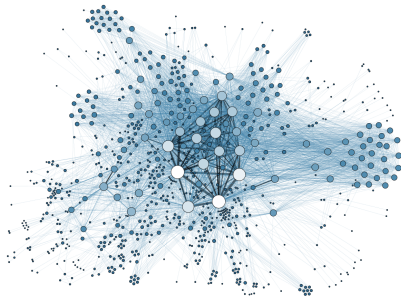
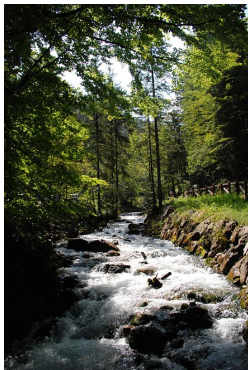
Electricity dataset (left), image [1] (right)

**Applications:** IoT, energy/traffic and demand prediction, monitoring and tracking, event and fraud detection, click/web logs, finance, reinforcement learning, . . . .

# Requirements

To deploy a model in the data stream setting, we require:

- 1 Prediction/action done **immediately** ( $\hat{y}_t = h_t(x_t)$  at time  $t$ )
- 2 Computational time spent per instance **must be less than the rate of arrival**



# Streaming Classification

Supervised ML models are often studied in the context of streams.

## Common assumptions found in the literature

- 1 Speed and size of stream implies instance-**incremental learning** (*at most* a single look at each data point)
- 2 The true label of data points become available (providing a stream of **training examples**)
- 3 **No temporal dependence**
- 4 **Concept drift** will occur

---

<sup>1</sup>e.g., predicting the weather – true label comes the next day

# Streaming Classification

Supervised ML models are often studied in the context of streams.

## Common assumptions found in the literature

- 1 Speed and size of stream implies instance-**incremental learning** (*at most* a single look at each data point)
- 2 The true label of data points become available (providing a stream of **training examples**)
- 3 **No temporal dependence**
- 4 **Concept drift** will occur

Some observations:

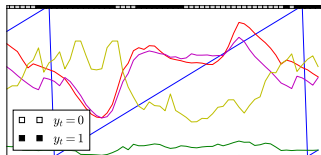
- Assumption 1 is unnecessary
- Assumption 2: Where do true labels from?
  - A human – then contradicts 1. (in most cases)
  - The future<sup>1</sup> – then contradicts 3. – it is a time series
- Assumptions 3 and 4 are contradictory

---

<sup>1</sup>e.g., predicting the weather – true label comes the next day

# Data Streams as Time Series

Benchmark datasets often look like time series:

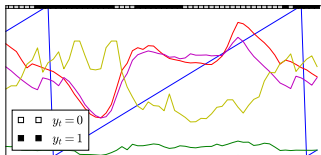


Prediction of Electricity demand



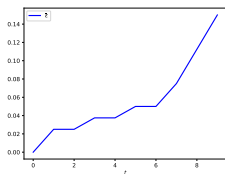
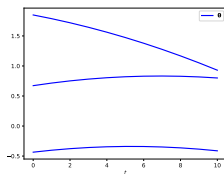
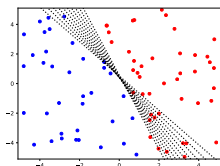
# Data Streams as Time Series

Benchmark datasets often look like time series:

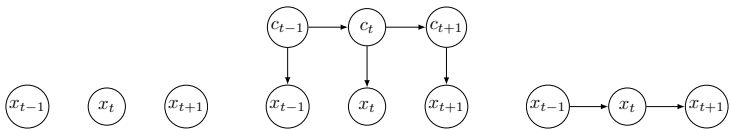


Prediction of Electricity demand

Even when points sampled iid wrt current concept, a time series forms in the coefficients, and/or in the error signal:



Concept drift  $\Rightarrow$  temporal dependence:



Time series tasks:

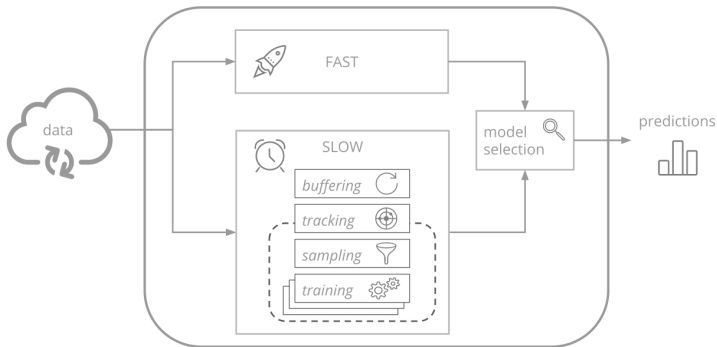
- Filtering
- Forecasting
- State labelling/change point detection
- Event/anomaly detection
- ...

(no supervised streaming classification!)

A **data stream** is a **time series** with constraints (prediction required *now*, update faster than rate of arrival).

# Fast and Slow Learning

- A framework for **Fast and Slow** learning
- Invest in higher level (slow) processes
- Batch and stream learning need not be mutually exclusive
- Time series methods, weakly labeled and unlabeled data
- Awareness of multi-input *multi-output* setting



Built into Scikit MultiFlow framework: <https://scikit-multiflow.github.io/>

# Outline

- 1 Data Streams as Time Series
- 2 Chaining Methods for Multi-Output Learning
- 3 Applications of Chaining in Data Streams

# Multi-label Learning

Input, e.g.,

$\mathbf{x} =$



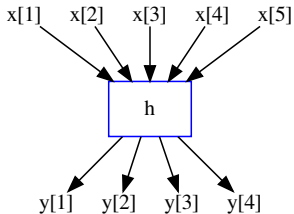
Prediction/output, e.g.,

$$\hat{\mathbf{y}} = [1, 0, 1, 0, 0] \Leftrightarrow \{\text{beach, foliage}\}$$

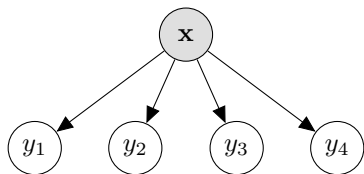
i.e., **multiple** outputs *per instance*.

Multi-label Problem  $[Y_1, \dots, Y_L] \in \{0, 1\}^L$

$X_1$	$X_2$	$X_3$	$X_4$	$X_5$	$Y_1$	$Y_2$	$Y_3$	$Y_4$
1	0.1	3	A	NO	0	1	1	0
0	0.9	1	C	YES	1	0	0	0
0	0.0	1	A	NO	0	1	0	0
1	0.8	2	B	YES	1	0	0	1
1	0.0	2	B	YES	0	0	0	1
0	0.0	3	A	YES	?	?	?	?



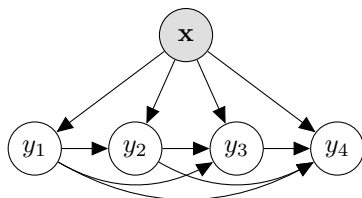
## Why not Independent Classifiers?



If we model labels together, we can achieve

- Better **predictive performance**
- Better **computational performance**
- Interpret **relationships** among labels (i.e., interpretability)
- Approach **structured-output prediction** tasks

# Classifier Chains



- Predictions **cascade along a chain** (as additional features)
- Has a probabilistic interpretation:

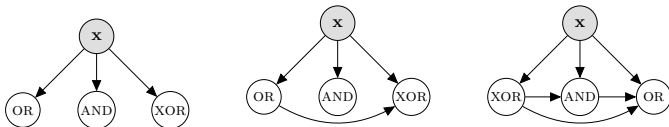
$$\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y} \in \{0,1\}^L} P(y_1|\mathbf{x}) \prod_{j=2}^L P(y_j|\mathbf{x}, y_1, \dots, y_{j-1})$$

- Inference becomes a **search** (for best  $\hat{\mathbf{y}}$ , in  $\{0,1\}^L$  space); e.g., greedy, Monte Carlo search,  $\epsilon$ -greedy, beam search.



# Ordering/Structuring the Labels

- Existing **hierarchy**? **May not be useful**
  - Only models positive dependence (if human-defined)
  - No guarantee of suitability for chosen classifiers
- Based on **label dependence**? **It depends** (on classifiers, inference, ...); consider:



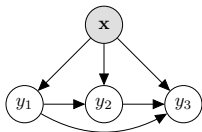
Metric	(left)	(middle)	(right)
Hamming score	0.83	1.00	0.83
Exact match	0.50	1.00	0.50

Logistic regression at each node  $h_j$ , greedy inference

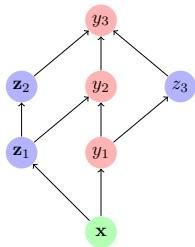
- Hill-climbing** in the label-structure space: Slow(!), but
  - Many local maxima (easy to reach) – i.e., **it works!**
  - Can make use of sub-optimal models that were trialled

# Classifier Chains: Why it Works

As a **probabilistic graphical model**:

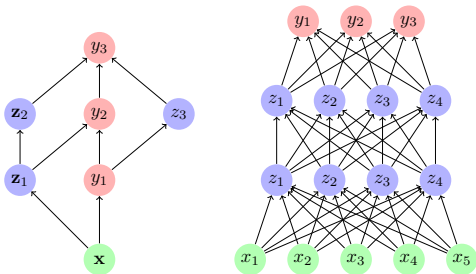


vs as a **neural network** ( $z$  delay nodes simply carry forward value):



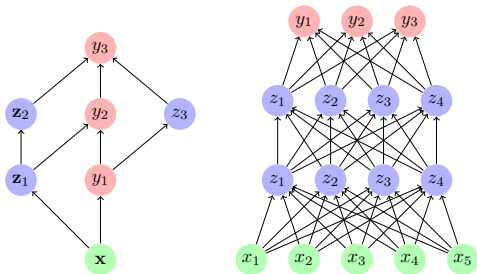
it's **deep in the label space!**

## Advantages vs standard neural network?



- Just apply 'off-the-shelf' [deep] neural net?
  - Dependence is modelled via the hidden layer(s)
  - Well-established, popular, competitive
- But with classifier chains:
  - The 'hidden' nodes come 'for free' (they're not hidden): faster training, less data required
  - A form of **transfer learning**

## Advantages vs standard neural network?

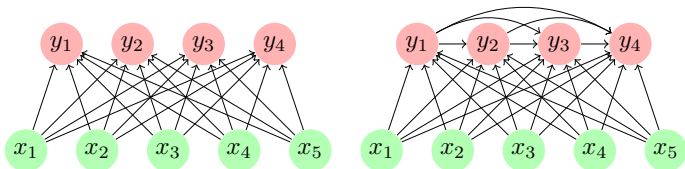


- Just apply 'off-the-shelf' [deep] neural net?
  - Dependence is modelled via the hidden layer(s)
  - Well-established, popular, competitive
- But with classifier chains:
  - The 'hidden' nodes come 'for free' (they're not hidden): faster training, less data required
  - A form of **transfer learning**

Observation: a bad/outdated prediction does not mean a bad representation!

# Multi-Output Regression

$X_1$	$X_2$	$X_3$	$X_4$	$X_5$	$Y_1$	$Y_2$	$Y_3$
1	0.1	3	A	NO	37.00	25	0.88
0	0.9	1	C	YES	-22.88	22	0.22
0	0.0	1	A	NO	19.21	12	0.25
1	0.8	2	B	YES	88.23	11	0.77
1	0.0	2	B	YES	?	?	?

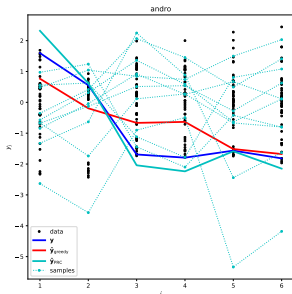
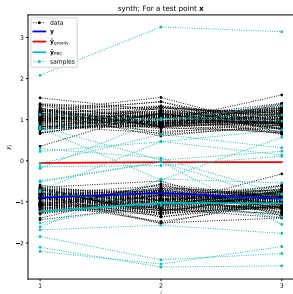


- Individual regressors – **directly applicable**.
- Chains
  - greedy inference – directly applicable, **but may be pointless!**
  - with probabilistic inference – **not tractable**, but we can sample *if we have  $p(y_j | \mathbf{x}, y_1, \dots, y_{j-1})$* .

# Regressor Chains

Results of chains under MSE (mean squared error) no better than using individual models / not interesting, *unless*

- Predictions provide an improved (**non-linear**) representation.
- **Non-isotropic** (state space models; where  $x_j$  is seen at 'time'  $j$ )
- We are interested in **interpretation**/explainability, e.g.,
  - Anomaly detection
  - Missing-value imputation
- New label concepts arrive later (we can **transfer learning**), make computational **time savings**.

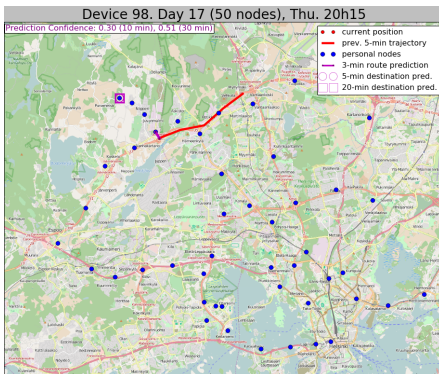


# Outline

- 1 Data Streams as Time Series
- 2 Chaining Methods for Multi-Output Learning
- 3 Applications of Chaining in Data Streams**

# Route Forecasting

- Create 'personal nodes' for a traveller
- Model and predict routes using classifier chains
- An advantage with relatively little training data and vs other methods (e.g., HMM, RNN)

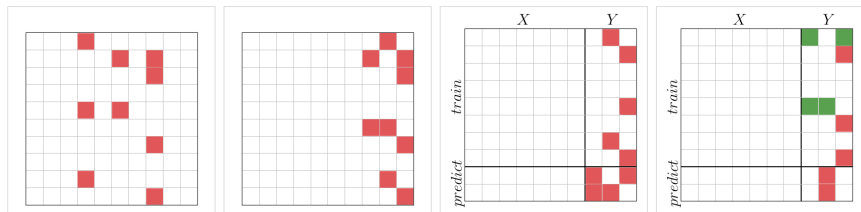


Personal nodes of a traveller and a predicted trajectory



# Missing-Value Imputation

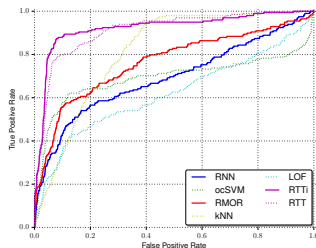
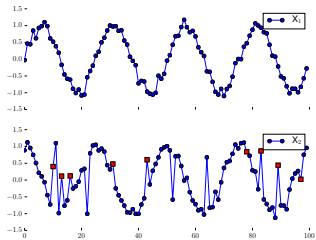
- Some values in the stream are missing!
- Turn the stream into multi-output samples, train, and predict (*impute*) missing values.
- Related to tasks in recommender systems



A set/stream of data transformed into a multi-output prediction problem.

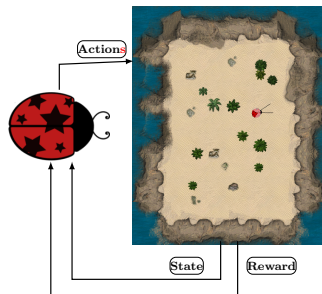
# Anomaly Detection and Interpretation

- Create 'random threads' (classifier/regressor chain cascades) through feature space and time (window)
- Monitor error spaces for anomalies
- Generate likely paths over the 'gap' (expand the number of samples if necessary)
- Impute this (treat it as a missing value) prior to using as a training example



# Continual Learning

In reinforcement learning,



- Reward signal is sparse
- Self-train on own **surrogate reward**, then use it as a feature.
- Recall: Incorrect predictions are not useless representations
- i.e., build up representation; transfer learning.

# Summary

- 1 Data Streams as Time Series
- 2 Chaining Methods for Multi-Output Learning
- 3 Applications of Chaining in Data Streams

# Unified Streaming/Batch Learning and Explainable Multi-output Prediction

Jesse Read



<http://www.lix.polytechnique.fr/~jread/>  
<http://www.lix.polytechnique.fr/dascim/>