# Online machine learning

Chris | Dong | Edvards | Hasan | Yang

## Introduction

- Data becomes available sequentially
- New data incrementally updates the model
- Useful when:
  - The entire dataset it too large
  - Patterns emerge dynamically
  - Data is generated over time

## Example applications

- Prediction from expert advice
  - Decider tries to perform as well as experts in hindsight
- Online spam filtering
  - Learning a binary classifier
- Online shortest paths in graph
  - Decider chooses the path
  - Adversary chooses the cost
- Portfolio selection
  - Decider chooses distribution of wealth over assets
  - Adversary chooses market returns
  - Decider learns to rebalance portfolio

## Basic concepts

- The framework is game-theoretic and adversarial
- For each iteration
  a. The decider makes a choice
  b. A convex cost function is revealed
  c. The decider incurs a cost
  d. The decider make a new choice to minimise regret
- Regret is the difference between the incurred cost and the cost of the best decision in hindsight

$$\text{regret} = \sum_{t=1}^{T} f_t(\mathbf{x}_t) - \min_{\mathbf{x} \in \mathcal{K}} \sum_{t=1}^{T} f_t(\mathbf{x})$$

## Convex optimisation

- We seek to minimise a continuous convex function over a convex subset of Euclidean space
- Gradient descent (GD) is the simplest and oldest optimisation method
- GD lays the foundation for more efficient and forthcoming algorithms

# First and second order methods

- Online gradient descent
  - Step in the direction of the gradient of the previous cost
  - If the new point is outside the underlying convex set, project it back within.
  - The regret is sublinear
  - But projection is burdensome
- Online Newton step
  - Approximates second derivative
  - Requires less iterations
  - But each step is costly
- Both algorithms require projection back into the convex set if they step out
- Projection is "expensive"

# Regularisation

- Follow the leader (FTL)
  - At any point in time, use the optimal decision in hindsight
  - Simple strategy
  - Regret is linear in iterations
  - Very unstable, changing decision too often
- Regularised FTL (RFTL)
  - Adds a regularisation function
  - Gives asymptotically optimal regret bounds
  - Stabilises the prediction

# Optimal regularisation

- We assume the regulariser is a strongly convex function, but which one?
  - It should depend on the decision set and cost function
- Adaptive subgradient method (AdaGrad)
  - Learns the optimal regulariser in hindsight online!

1: Input: parameters $\eta, \mathbf{x}_1 \in \mathcal{K}$.
2: Initialize: $S_0 = G_0 = \mathbf{0}$,
3: **for** $t = 1$ to $T$ **do**
4:     Predict $\mathbf{x}_t$, suffer loss $f_t(\mathbf{x}_t)$.
5:     Update:
$$S_t = S_{t-1} + \nabla_t \nabla_t^\top, \ G_t = S_t^{1/2}$$
$$\mathbf{y}_{t+1} = \mathbf{x}_t - \eta G_t^{-1} \nabla_t$$
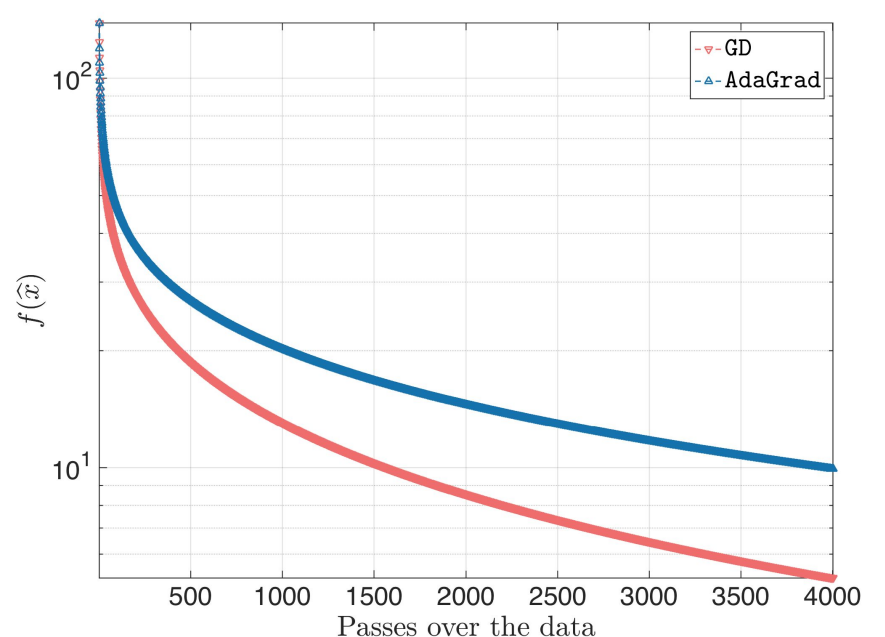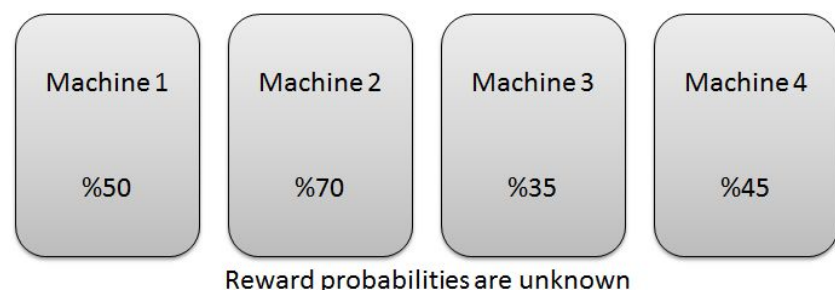$$\mathbf{x}_{t+1} = \arg\min_{\mathbf{x} \in \mathcal{K}} \|\mathbf{y}_{t+1} - \mathbf{x}\|_{G_t}^2$$
6: **end for**

## Online decision-making Bandit & Reinforcement learning



| Machine 1 | Machine 2 | Machine 3 | Machine 4 |
|-----------|-----------|-----------|-----------|
| %50 | %70 | %35 | %45 |

Reward probabilities are unknown

- At each step $t=1,2,...,T$, a decision-maker
  - Observes the state,
  - Chooses an action from a given action set $A$,
  - Receives reward.
- Goal is to maximize the collected rewards.
- *Regret for decision-making:*

$$Regret_{\pi,T} := \mathbb{E}^{\pi^*}\left\{\sum_{t=1}^{T} r_t\right\} - \mathbb{E}^{\pi}\left\{\sum_{t=1}^{T} r_t\right\}$$

- *Regret for RL in MDPs:*

$$Regret_{\mathbb{A},T}(s_1) := Tg^{\star}(s_1) - \sum_{t=1}^{T} r(s_t, a_t)$$

- *Can't be arbitrarily minimized due to some fundamental performance limits.*

A simple multi armed bandit algorithm:

- With some probability, explore the action space.
  - Use the feedback to construct an estimate of the actions' losses.
- Otherwise, use the estimates to select the optimum choice.
  - Suppose the estimates are the true historical costs.

## Projection-free algorithms

- In many computational and learning scenarios the main bottleneck of optimization, both online and offline, is the computation of projections onto the underlying decision set.
- The conditional gradient(CG) method, or Frank-Wolfe algorithm, is a simple algorithm for minimizing a smooth convex function f over a convex set.

1: **Input:** step sizes $\{\eta_t \in (0,1), \ t \in [T]\}$, initial point $\mathbf{x}_1 \in \mathcal{K}$
2: **for** $t = 1$ to $T$ **do**
3:     $\mathbf{v}_t \leftarrow \arg\min_{\mathbf{x} \in \mathcal{K}} \left\{ \mathbf{x}^\top \nabla f(\mathbf{x}_t) \right\}$.
4:     $\mathbf{x}_{t+1} \leftarrow \mathbf{x}_t + \eta_t (\mathbf{v}_t - \mathbf{x}_t)$.
5: **end for**

- Matrix completion problem:

$$\min_{X \in \mathbb{R}^{n \times m}} \frac{1}{2}\|X - M\|_{OB}^2$$

$$\text{s.t.} \quad \text{rank}(X) \leq k.$$
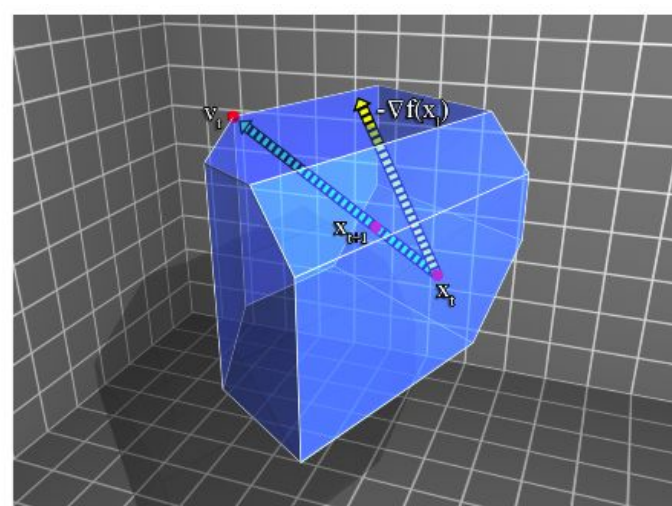
- CG for matrix completion

1: Let $X^1$ be an arbitrary matrix of trace $k$ in $\mathcal{K}$.
2: **for** $t = 1$ to $T$ **do**
3:     $\mathbf{v}_t = \sqrt{k} \cdot v_{\max}(-\nabla_t)$.
4:     $X^{t+1} = X^t + \eta_t(\mathbf{v}_t \mathbf{v}_t^\top - X^t)$ for $\eta_t \in (0,1)$.
5: **end for**
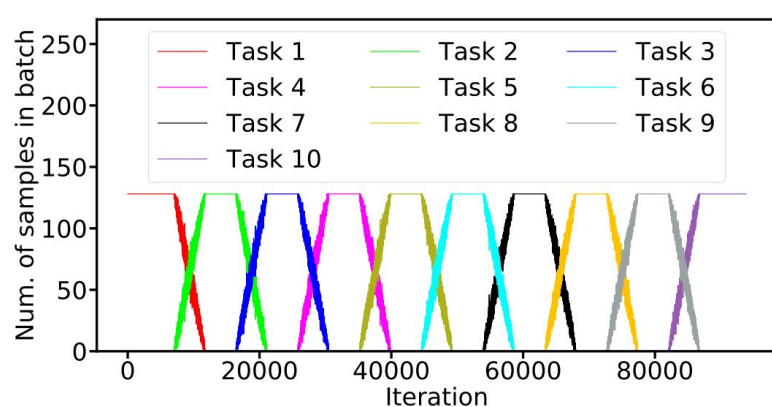
In the CG method, the update to the iterate $\mathbf{x}_t$ may not be in the direction of the gradient, as $\mathbf{v}_t$ is the result of linear optimization procedure.

# Online Bayes learning

- Common continuous learning issues:
    - Data distribution is changing while learning (Catastrophic forgetting problem)
    - Multiple related tasks while no clear boundaries
    - Asynchronous data arrival
    - Reliability of gradient



- Bayes rule: encode past information into posterior, which is used as prior for future predictions

- Practical solutions for posteriors estimation:
    - Variational methods
    - Monte Carlo methods
    - Laplace/Mean-field approximations/ADF/EP
- Encode estimation belief into learning rate/speed:
    - Bayesian Gradient Descent

$$p\left(\boldsymbol{\theta}|D_n\right) = \frac{p\left(D_n|\boldsymbol{\theta}\right)p\left(\boldsymbol{\theta}|D_{n-1}\right)}{p\left(D_n\right)}$$

# Conclusion

Online machine learning is useful in many different applications where data becomes available over time.

Some examples relating to our research:

- Adaptive control (e.g., under changing dynamics)
- Training over large datasets (e.g., in imitation learning)
- Reinforcement learning over a network of agents