# ML 540 Tutorial

Reza Babanezhad

rezababa@cs.ubc.ca

UBC

# Outline

- Linear Programming
  - Some example
- MAP estimation with different distributions
- Fun with ML

# Linear Programming

- Linear programming has two main component:
  - Linear Objective function
  - Linear constraints could including both equalities and inequalities

$$let \quad a_{ji}, c_i, x_i, b_j \in \mathbb{R} \quad, x^T = (x_1, \ldots, x_n) \quad i \in [1..n]$$
$$j \in [1..m]$$

$$objective \ function: \quad \min_x \sum_{i=1}^{n} c_i x_i$$

$$s.t. \quad \sum_{i=1}^{n} a_{ji} x_i \leq b_j \quad \forall j \in [1..m]$$

$$x_i \geq 0$$

- In matrix form

Let $A \in \mathbb{R}^{m \times n}, b \in \mathbb{R}^m, x, c \in \mathbb{R}^n$

then

$$\min \quad c^T x$$

$$\text{s.t.} \quad Ax \leq b$$

$$x \geq 0 \quad \longleftarrow \text{ This is extra}$$

constraint on $x$!

- To deal with extra constraint on x
  - Divide x into positive and negative parts

$$x = x^+ - x^-$$

$$x_i^+ = \max(x_i, 0) \geq 0, \quad x_i^- = \max(-x_i, 0) \geq 0$$

objective : $\min \ c^T x^+ - c^T x^-$

s.t. $\quad A x^+ - A x^- \leq b$

$$x^+ \geq 0$$

$$x^- \geq 0$$

- A little geometry
  - Hyperplane : C'x=b is a hyperplane    $C, x \in \mathbb{R}^n$
    - For example 2x+3y+5z=4 is a hyperplane in 3D space.
    - Hyperplane is a convex set- if we connect two points of the set, the entire line is still in the set
  - Each hyperplane divides the space into 2 half spaces: C'x <= b or C'x > b
    - Half space is a convex set
  - Intersection of two convex set is still convex. (Why?)
  - Polytope: intersection of some half spaces and hyperplane:    $Ax \leq b$
    - Polytope is a convex set
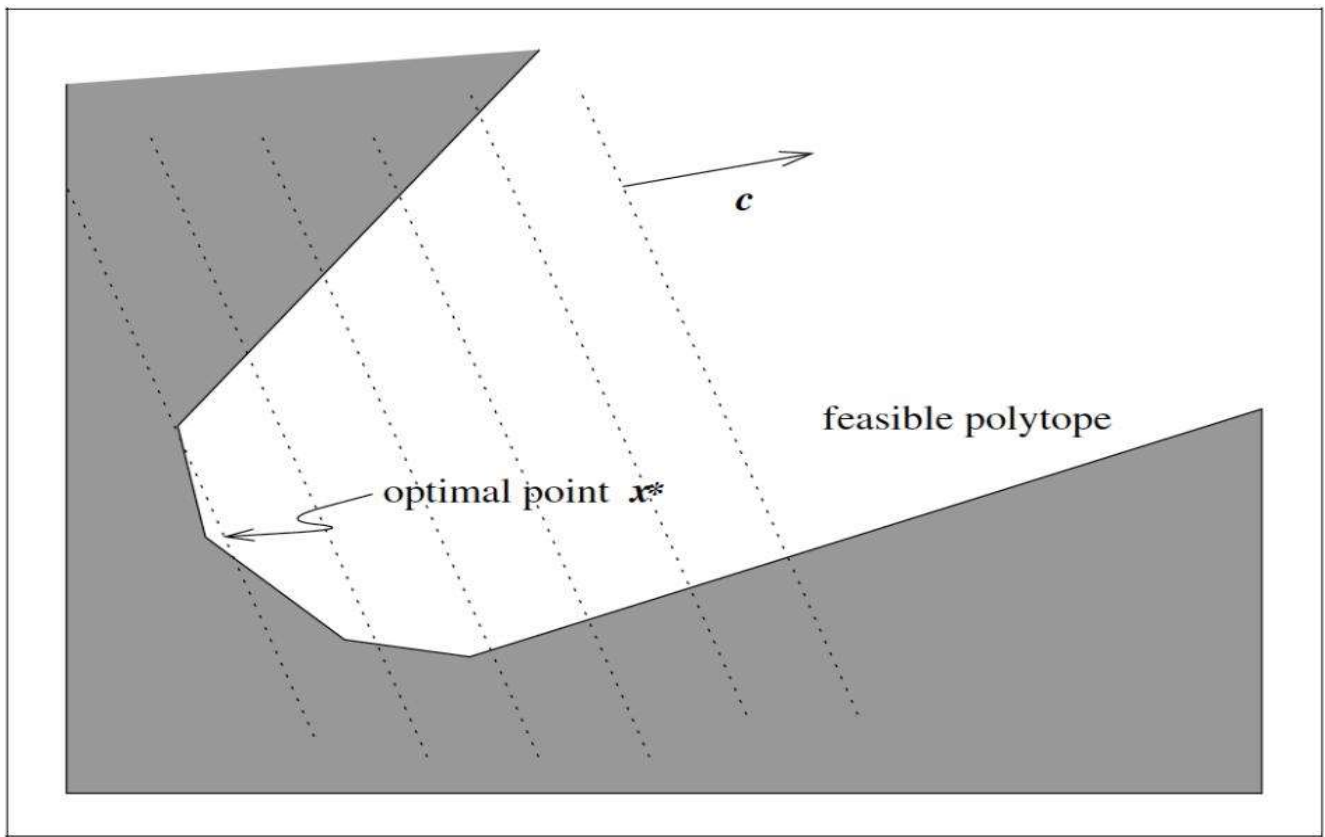
$A \in \mathbb{R}^{m \times n}$

$b \in \mathbb{R}^m$

- How can we interpret LP with geometrical objects?
- We can set C'x=z in our LP so
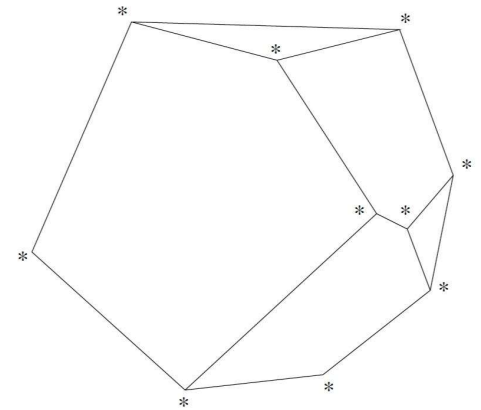- We want to find a hyperplane which its intersection with a polytope is minimum among all other hyperplanes

$$\min \quad z$$
$$\text{s.t} \quad \begin{cases} Ax \leq b \\ x \geq 0 \\ c^T x = z \end{cases}$$

half-space

Polytope

hyperplane

$c$

feasible polytope

optimal point $x^*$

- How to solve LP?
  - Simplex Method
    - All feasible solutions are vertices of the feasible polytope
    - Cost in worst case: exponential
  - Interior point methods
    -  formulate as non-linear problem
    - Polynomial time in worst case e.g. O(n^4L) for ellipsoid method
  - MATLAB uses Interior-Point-Legacy Algorithm

- MATLAB command for LP

$$\min_{x} \; c^T x$$

$$\text{s.t.} \; A x \le b$$

$$A_{eq} x = b_{eq}$$

$$lb \le x \le ub$$

$$x = linprog(c, A, b)$$

$$x = linprog(c, A, b, A_{eq}, b_{eq})$$

$$x = linprog(c, A, b, A_{eq}, b_{eq}, lb, ub)$$

e.g.1. $\arg\min\limits_{x} |a_i^T x + b|$   $a_i \in \mathbb{R}^d$ , $x \in \mathbb{R}^d$, $b \in \mathbb{R}$

1_ Introduce new variable $r \in \mathbb{R}$

$$r \geq \{ax + b, -ax - b\}$$

2_ rewrite the problem

$$\min\limits_{r} r$$
$$ax - r \leq -b$$
$$-r - ax \leq b$$

3- find approperiate prameters for lin prog : $C, A, b$?

$$\min_{r, x} \quad r \quad \Rightarrow \quad r = \vec{0} \cdot x + 1 \times r \Rightarrow C^T = [\underbrace{0 \cdots 0}_{d} \; 1]_{1 \times (d+1)}$$

$$a x - r \leq -b$$
$$-r - a x \leq b$$

$$\Rightarrow \underbrace{\begin{bmatrix} a & -1 \\ -a & -1 \end{bmatrix}}_{\underset{A}{\underbrace{\phantom{aaaa}} 2 \times (d+1)}} \underbrace{\begin{bmatrix} x \\ r \end{bmatrix}}_{d+1 \times 1} \leq \underbrace{\begin{bmatrix} -b \\ b \end{bmatrix}}_{\underset{b}{\underbrace{\phantom{aa}}} 2 \times 1}$$

**e.g. 2 :**  $\min\limits_{x} |a_1 x + b| + |a_2 x + b| + |a_3 x + b|$

$$a_i, x \in \mathbb{R}^d , \quad b \in \mathbb{R}$$

1: Introduce $r_i$ corrosponding to each $|a_i x + b|$ s.t. $r_i \in \mathbb{R}$

and $r_i \geq \max\{a_i x + b, -a_i x - b\}$

2- rewrite LP

$$\min\limits_{r_i, x} r_1 + r_2 + r_3$$

$$-r_i - a_i x \leq b_i$$
$$-r_i + a_i x \leq -b_i \qquad \forall i$$

3 - find $c, A, b$

$$r_1 + r_2 + r_3 = \vec{0}^T x + 1 \cdot \vec{r} \implies c^T = [\vec{0}^T_{1\times d} \quad \vec{1}^T_{1\times 3}]$$

$-r_i - a_i^T x \le b_i \quad \forall i$

$-r_i + a_i^T x \le -b_i \quad \forall i$
$\implies$

$$\underbrace{\begin{bmatrix} a_1^T & -1 & 0 & 0 \\ a_2^T & 0 & -1 & 0 \\ a_3^T & 0 & 0 & -1 \\ -a_1^T & -1 & 0 & 0 \\ -a_2^T & 0 & -1 & 0 \\ -a_3^T & 0 & 0 & -1 \end{bmatrix}}_{A} \underbrace{\begin{bmatrix} x \\ r_1 \\ r_2 \\ r_3 \end{bmatrix}}_{(d+3)\times 1} \le \underbrace{\begin{bmatrix} -b_1 \\ -b_2 \\ -b_3 \\ b_1 \\ b_2 \\ b_3 \end{bmatrix}}_{b \quad 6\times 1}$$

$6 \times (d+3)$

- let $I_n$ to be $n \times n$ identity matrix

let $\begin{bmatrix} a_1^T \\ a_2^T \\ a_3^T \end{bmatrix} = \alpha$ $\qquad \beta = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix}$

then $A = \begin{bmatrix} \alpha & -I \\ -\alpha & -I \end{bmatrix}$, $b = \begin{bmatrix} -\beta \\ \beta \end{bmatrix}$

In Matlab: $A = [\alpha \quad -eye(3); \quad -\alpha \quad -eye(3)]$, $b = [-\beta; \beta]$

e.g.: $\text{argmin}_x \text{Max}_{\hat{a}} \{|a_i x + b_i|\}$ $i \in [1 .. n]$

Let $r$ be the minimum of the LP so for all $i$

we have: $r \geq |a_i x + b_i|$ since $r = \max_i \{|a_i x + b_i|\}$

So the LP would be

$$\min r$$
$$\text{s.t.} \quad r \geq a_i x + b_i$$
$$r \geq -a_i x - b_i \quad \forall i$$

and from previous examples we can solve this one too!

# MAP Estimation

Let $y_i \sim N(\mu_i, 1)$ & $\mu_i \sim N(0, \eta_i)$   $i \in [1..n]$

MAP for $\mu$?

For MAP we need likelihood $\times$ prior

So $\prod_{i=1}^{n} N(\mu_i, 1) \cdot N(0, \eta_i) = \mathcal{L}(\mu)$

now taking the Log $\left( \log \prod = \sum \log \right)$

$$\ell(\mu) = \sum_{i=1}^{n} \log N(\mu_i, 1) + \log N(0, \eta_i)$$

$$\log \mathcal{N}(\mu_i, 1) = -\frac{1}{2}(y_i - \mu_i)^2$$

$$\log \mathcal{N}(0, \eta_i) = -\frac{1}{2}\eta_i^{-1}\mu_i^2$$

So $\ell(\mu) = \sum_{i=1}^{n} -\frac{1}{2}(y_i - \mu_i)^2 - \frac{1}{2}\eta_i^{-1}\mu_i^2$

if $y^T = (y_1, \dots, y_n)$, $\mu^T = (\mu_1, \dots, \mu_n)$ and $\eta^{-1} = \begin{bmatrix} \eta_1^{-1} & & 0 \\ & \ddots & \\ 0 & & \eta_n^{-1} \end{bmatrix}_{n \times n}$

$$\ell(\mu) = \frac{-1}{2}\|y - \mu\|^2 - \frac{1}{2}\mu^T \eta^{-1} \mu$$

Let $y_i \sim \mathcal{L}(\mu_i, 1) = \frac{1}{2} \exp(-|y_i - \mu_i|)$, $\mu_i \sim \mathcal{N}(0, \eta)$

MAP for $\mu$?

likelihood $\times$ prior $= \prod_{i=1}^{n} \mathcal{L}(\mu_i, 1) \times \mathcal{N}(0, \eta) = L(\mu)$

$-\log L(\mu) = \sum_{i=1}^{n} |y_i - \mu_i| + \frac{\eta^{-1}}{2} \mu_i^2 + c \longrightarrow$ constant

$\Rightarrow$ MAP : $\underset{\mu}{\arg\min} \; \|y - \mu\| + \frac{\eta^{-1}}{2} \|\mu\|_2^2$

Consider $y_i$ comes from a Student t with $v$ degree of freedom. And the mean of r.v. $y_i$ is $w^T x_i$; & $W$ is the parameter of the model and $w \sim N(0, \eta I)$. Student t distribution function with $0$ mean i.e. $E[\Theta] = 0$ is:

$$P(\Theta | v) = c \left(1 + \frac{\Theta^2}{v}\right)^{-\frac{v+1}{2}}$$

Find MAP for $w$?

$E[y_i] = W^T x_i \neq 0$ so to use Student t PDF we set

$\Theta_i = y_i - W^T x_i$ and let scale parameter to be 1

$$\text{So} \quad P(\Theta_i \mid v) = C \left( 1 + \frac{|y_i - w^T x_i|^2}{v} \right)^{-\frac{v+1}{2}}$$

$$\text{likelihood} \times \text{prior} = \prod_{i=1}^{n} C \left( 1 + \frac{|y_i - w^T x_i|^2}{v} \right)^{-\frac{v+1}{2}} \times N(0, \eta) = L(w)$$

$$-\log L(w) = \sum \frac{v+1}{2} \log \left( 1 + \frac{|y_i - w^T x_i|^2}{v} \right) + \frac{\eta^{-1}}{2} \|w\|_2^2$$

# Fun Time

- Assume a start up company hired you as a data scientist to design and implement a recommendation system for them. This company like Youtube allows its users to upload videos. So your job is to build a simple recommendation system considering the list of watched videos and given number of views by each user, recommend new videos to users.

# Set-based approach

- Ignore watching number!
- For each video build a set of users
- Find the intersections of these user sets for all videos
- Recommend the video to the users who watched another video whose intersection with this one is big enough.
- What do you think about this model? (population effect?)

# Cosine based

- Assume each video as a vector in user space.
- Each dimension shows the number times that video is watched by corresponding user.
- Find the cosine between vectors.
- Rank related videos based on cosine for each videos.
- What is the problem?
  - Ignoring the overall activity of each user- more selective users versus listening to everything

# TF-IDF based

- Treat each user as term
- Treat each video as document
- TF: how many times each user watched a video
- IDF: how many times a user watched videos on the site
- Now build a vector of TF-IDF weight for each video
- Find the cosine ☺