

FCM application
Experience series
 (x_i, y_i)
 $x_i \in \mathbb{R}^d$
 $y_i \in \{-1, 1\}$
 \rightarrow $y_i = f(x_i)$

Expertise
 $y = f(x)$
 $f: \mathbb{R}^d \rightarrow \{-1, 1\}$
Experience $\xrightarrow{\text{ML}}$ Expertise
 f is designed such that $y_i = f(x_i)$
 $\forall i=1:N$

course is all about finding f to get $y = f(x)$

ML algos: issues
scalability
stability
complexity
for finding
evalues of $A^{d \times d}$
Not scalable
Spam account + account or account
should still recognize
learning sup - labels available
unsup - unavailable labels

learning algos: input \Rightarrow data
issues in ML

One major issue is datum size can be addressed through

- Dimensionality reduction (DR)
- feature extraction or efficient representation of data

Several DR methods exist: SVD/PCA, manifold learning

MBML

Other issues deal w/ designing learning algos, validating such algos etc
Another example

$x_i \in \mathbb{R}^d \rightarrow$ parameters on
 $i=1:N$ loan seeker

$y_i \in \{-1, 1\} \rightarrow$ whether loan seeker is likely to default or not.

$x_i \in \mathbb{R}^d \rightarrow z_i \in \mathbb{R}^n$ need projection to lower dim space

when projecting to lower dim space we want discriminating features are retained.

Ideally, want $\|x_i - z_j\|_2 = \|z_i - z_j\|_2$

since, equality is not possible

at least $\|x_i - z_j\|_2 \approx \|z_i - z_j\|_2$

$$z_i = g(x_i)$$

vector in lower dim space \rightarrow higher dim vector
 $\rightarrow g$ has to be built

We may have to reproject \rightarrow back to higher dim.

$$x_i \in \mathbb{R}^d \rightarrow z_i \in \mathbb{R}^n \text{ (need)} \rightarrow \tilde{x}_i$$

patient in remote region \rightarrow lower dim for transmission \rightarrow compression \rightarrow reconstruction \rightarrow Expert needs higher dim again

$$\text{Expectation } \rightarrow x_i \approx \tilde{x}_i$$

If too much info is lost, we don't face issues

we want small value for n \rightarrow good reconstruction. \rightarrow strike a balance

So in $z_i = g(x_i)$, g can't be arbitrary

PCA

method of random projection

linear (simple)

non-linear manifold learning (can get small n)

Linear vs non linear

Simple
can't get by small $n \times n$ (lemma)

can get small $n \times n$
Creators of manifold learning

projected faces to get moods on 2D plane

Summary

$$(x_i, y_i)_{i=1}^n, x_i \in \mathbb{R}^d$$

if d is big

Dimensionality reduction

linear nonlinear

- PCA
- random projection
- Fractals
- manifold learning

$$(x_i, y_i)_{i=1}^n \quad \text{small dim.}$$

Learning

x_i

supervised

$$(x_i, y_i)$$

label

unsupervised

K-means

Sparcely based clustering

SVM

Neural nets

SVD

$$A_{m \times n} \quad \text{rank}(A) = r \leq \min(m, n)$$

$$\text{SVD: } A = U \Sigma V^T \quad \begin{matrix} m \times m \\ n \times n \end{matrix} \quad (\text{both being unitary})$$

unitary) and a "diagonal" matrix $\Sigma_{m \times n}$ such that $A = U \Sigma V^T$

Proof of $A = U \Sigma V^T$

$$(A^T A) \quad \text{symmetric}$$

Its eigenvalues must be ≥ 0

$$(A^T A) v_i = \lambda_i^2 v_i$$

$$\lambda_i^2 \geq 0 \quad i > r \quad (\text{where } r \text{ is rank})$$

Define $u_i \in \mathbb{R}^m$ s.t.

$$u_i = \frac{1}{\sqrt{\lambda_i}} A v_i \quad i=1 : r$$

$$u_i^T u_j = \frac{1}{\sqrt{\lambda_i} \sqrt{\lambda_j}} (A v_i)^T (A v_j)$$

$$= \frac{1}{\sqrt{\lambda_i} \sqrt{\lambda_j}} v_i^T (A^T A) v_j$$

$$= \frac{1}{\sqrt{\lambda_i} \sqrt{\lambda_j}} v_i^T \cancel{\lambda_i^2} \cancel{v_i} (x_j - v_j)$$

$$u_i^T u_j = \frac{v_i^T v_j}{\sqrt{\lambda_i}}$$

$$= \frac{v_i^T v_j}{\sqrt{\lambda_i}} \quad \begin{cases} 1 & i=j \\ 0 & i \neq j \end{cases}$$

u_1, u_2, \dots, u_r are LI

$\{u_1, u_2, \dots, u_r\}$ can be extended to a basis say $\{u_1, u_2, \dots, u_m\}$

$A V$

$$= [AV_1 \quad AV_2 \quad \dots \quad AV_n]$$

$$= [\lambda_1 u_1 \quad \lambda_2 u_2 \quad \dots \quad \lambda_r u_r \quad \cancel{0} \quad \cancel{0} \quad \dots]$$

$$= \begin{bmatrix} \lambda_1 & & & & \\ & \lambda_2 & & & \\ & & \ddots & & \\ & & & \lambda_r & \\ & & & & 0 \end{bmatrix} \begin{bmatrix} u_1 & u_2 & \dots & u_m \end{bmatrix} \begin{bmatrix} 1 & & & & \\ & 1 & & & \\ & & \ddots & & \\ & & & 1 & \\ & & & & 0 \end{bmatrix}$$

$$AV = UD$$

$$\text{or } A = UDV^T$$

Aug 04

MBM

$$f(x) : \mathbb{R}^n \rightarrow \mathbb{R}$$

Directional derivative

Suppose $d \in \mathbb{R}^n$

$$f'_d = \lim_{t \rightarrow 0} \frac{f(x+td) - f(x)}{t}$$

Suppose $d = e_i = (1, 0, 0, \dots, 0)$

$$\text{then } f'(x, e_i) = \frac{\partial f}{\partial x_i}$$

directional derivative along d

Suppose all partial derivatives exist for f , then gradient of f is defined by

$$\nabla f = \begin{pmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{pmatrix} \quad \text{gradient}$$

$$f'(x^*, d) = \nabla f(x^*) \cdot d$$

Suppose double partial derivatives exist for f (continuous), at all points, Hessian is defined by

$$\nabla^2 f(x) = \begin{pmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \vdots & & & \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{pmatrix}$$

Basic optimization theory

The following are some opt. problems

$$\min_{u, w} \sum_{i=1}^n \|x_i - u w x_i\|_2^2$$

reconstruction

compression

unconstrained optimization problem

$$\min_{u_i \in S} \sum_{i=1}^k \|x_i - u_i\|$$

$$\min \left(\frac{1}{2} \|w\|_2^2 \text{ such that } g(w, t) \geq 1 \right)$$

constrained optimization problem

Basic opt theory

$f: \mathbb{R}^n \rightarrow \mathbb{R}$ is ~~continuous~~

and $0 \neq d \in \mathbb{R}^n$

$$\text{if } \lim_{t \rightarrow 0} \frac{f(x_0 + td) - f(x_0)}{t} \text{ exists}$$

$$= f'(x_0, d)$$

(notation)

directional derivative of f at x_0 along direction d

If $d = e_i = (1, 0, 0, \dots, 0)$, then

$$f'(x_0, d) = \frac{\partial f}{\partial x_i}(x_0)$$

If $f'(x_0, e_i)$ exists $\forall i=1:n$,

$$\nabla f(x_0) = \begin{pmatrix} \frac{\partial f(x_0)}{\partial x_1} \\ \vdots \\ \frac{\partial f(x_0)}{\partial x_n} \end{pmatrix} \quad \text{gradient of } f \text{ at } x_0$$

$$f'(x_0, d) = \nabla f(x_0) \cdot d$$

Suppose f has continuous partial derivatives of second order

$$H = [\nabla^2 f(x_0)] = \begin{pmatrix} \frac{\partial^2 f(x_0)}{\partial x_1^2} & \frac{\partial^2 f(x_0)}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f(x_0)}{\partial x_1 \partial x_n} \\ \vdots & & & \\ \frac{\partial^2 f(x_0)}{\partial x_n \partial x_1} & \cdots & \frac{\partial^2 f(x_0)}{\partial x_n^2} \end{pmatrix}$$

(notation)

$$\begin{pmatrix} \frac{\partial^2 f(x_0)}{\partial x_1^2} & \cdots & \frac{\partial^2 f(x_0)}{\partial x_1 \partial x_n} \\ \vdots & & \\ \frac{\partial^2 f(x_0)}{\partial x_n \partial x_1} & \cdots & \frac{\partial^2 f(x_0)}{\partial x_n^2} \end{pmatrix}_{m \times n}$$

Hessian of f at x_0 .

First order approximation theorem

Suppose $f: \mathbb{R}^n \rightarrow \mathbb{R}$ has all continuous derivative of 2nd order and $x, y \in \mathbb{R}^n$.

Then $\exists z \in [x, y]$



$$f(y) = f(x) + (y-x)^T \nabla f(x) + \frac{1}{2}(y-x)^T \nabla^2 f(z)(y-x)$$

Second order approximation theorem

Suppose $f: \mathbb{R}^n \rightarrow \mathbb{R}$ has continuous partial derivatives of 2nd order

Then $f(y) =$

$$\begin{aligned} f(y) &= f(x) + (y-x)^T \nabla f(x) \\ &\quad + \frac{1}{2}(y-x)^T (\nabla^2 f(x))(y-x) \\ &\quad + o(\|y-x\|^2) \end{aligned}$$

↓ accounting for residual

little-o notation

$$f(n) = o(g(n))$$

$$g(a) \neq 0, \lim_{x \rightarrow a} \frac{f(x)}{g(x)} = 0$$

$$f_n = o(g_n), \lim_{n \rightarrow \infty} \frac{f_n}{g_n} = 0$$

Suppose $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is "sufficiently nice" (all partial derivatives must exist)

Theorem

If x^* is a local minimal point of f , then $\nabla f(x^*) = 0$

vector.

proof: $g(t) = f(x^* + te_i)$ for some $i \in \{1, \dots, n\}$

$$g(t) = f(x^* + te_i) \geq f(x^*) = g(0)$$

↑ t. such that

$x^* + te_i$ lies in the nbhd of x^*

$\Rightarrow 0$ is a local min pt of f or g

$$\therefore g'(0) = 0$$

$$g(t) = \frac{\partial f}{\partial x_i}(x^* + te_i)$$

$$g'(0) = \frac{\partial f}{\partial x_i}(x^*) = 0$$

$$\Rightarrow \nabla f(x^*) = 0$$

∴ If x^* is an optimal point
then $\nabla f(x^*) = 0$

How to know if a point x^* is maximum or minimum

we get from 2nd order approx theorem

Results

1) If x^* is a local minimum pt for f , then $(\nabla^2 f(x^*))$ is positive semi definite (PSD)

A is PSD if $x^T A x \geq 0 \forall x \in \mathbb{R}^n$

proof

from 2nd order approx theorem

$$f(y) = f(x) + (y-x)^T \nabla f(x)$$

$$+ \frac{1}{2}(y-x)^T (\nabla^2 f(x))(y-x)$$

$$+ o(\|y-x\|^2)$$

$$0 \leq f(y) - f(x^*) = \frac{1}{2}(y-x^*)^T \nabla^2 f(z)(y-x^*)$$

Aug 07

Recall:

$$\min_{x \in \mathbb{R}^n} f(x)$$

$$\mathbb{R} = \mathbb{R}^n$$

$$S \subseteq \mathbb{R}^n$$

unconstrained optimization problem

If f has a local min at $\underset{\text{max}}{x^*}$ then $\nabla f(x^*) = 0$

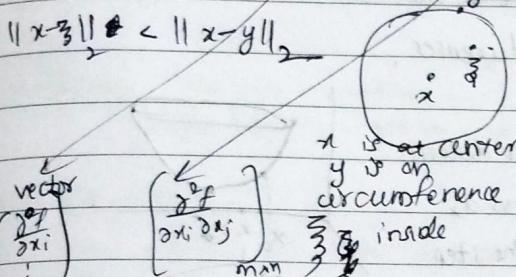
1st order approximation theorem

f is sufficiently nice

$$x, y \in \mathbb{R}^n \rightarrow$$

$\exists "z \in [x, y]"$ such that

$$f(y) = f(x) + (y-x)^T \nabla f(x) + \frac{1}{2}(y-x)^T \nabla^2 f(z)(y-x)$$



2nd order approximation theorem

$$f(y) = f(x) + (y-x)^T \nabla f(x) + \frac{1}{2}(y-x)^T \nabla^2 f(x) y(x) + o(\|y-x\|^2)$$

If f has a local min pt at x^* , then $\nabla^2 f(x^*) \geq 0$
positive semidefinite
(all eigen values ≥ 0)

$$f(y) - f(x) = \frac{1}{2}(y-x)^T \nabla^2 f(z)(y-x)$$

$$y \in B(x^*, r)$$

ball small

Tell y is close to x^* .

for all y close to x^*

If $\nabla^2 f(x)$ is PSD matrix, so

$$f(y) - f(x^*) \geq 0 \quad \forall y$$

x^* is local optima minima

If f has a local max at x^*

then $[\nabla^2 f(x)] \leq 0$.

If $[\nabla^2 f(x^*)] > 0$, then f has a strict min pt at x^* .

Saddle point

$$\nabla f(x^*) = 0$$

$[\nabla^2 f(x^*)]$ is neither PSD nor PND (positive negative definite)

Global min pt

$$f: \mathbb{R}^n \rightarrow \mathbb{R}$$

$$[\nabla^2 f(x)] \geq 0 \quad \forall x \in \mathbb{R}^n$$

The critical point x^* is a global min pt.

Example

$$f(x) = x^T Ax + 2b^T x + c, b \in \mathbb{R}^n, c \in \mathbb{R}$$

$$\text{gradient } \nabla f(x) = 2Ax + 2b$$

$$\text{Hessian } \nabla^2 f(x) = A$$

$$\nabla f(x) = 0 \Rightarrow Ax = -b$$

If $\nabla^2 f(x) = A \geq 0$ (PSD) $\forall x \in \mathbb{R}^n$

then each solution of ① is a min pt for \mathbb{R}^n

If $\nabla^2 f(x) = A \leq 0$ (PND) $\forall x \in \mathbb{R}^n$

then each solⁿ of ① is a global max pt for \mathbb{R}^n

If A is PD (strictly PD) (all eigen values are strictly > 0)

then A is invertible and

~~$x = -A^{-1}b$~~ is the critical point

If " $A > 0$ " $x = -A^{-1}b$ is the global min pt
 If " $A \leq 0$ " $x = -A^{-1}b$ is the global max pt

Theorem

Suppose $f(x) = x^T Ax + 2b^T x + c$, ($b \in \mathbb{R}^n$, $c \in \mathbb{R}$
 and A is symmetric)

Then the following are equivalent.

$$\textcircled{1} \quad (x^T x) \geq 0 \quad \forall x \in \mathbb{R}^n$$

$$\textcircled{2} \quad \begin{pmatrix} A & b \\ b^T & c \end{pmatrix} \geq 0 \quad \rightarrow (\text{is PSD})$$

Suppose \textcircled{2} holds.

$$(x) \begin{pmatrix} A & b \\ b^T & c \end{pmatrix} \begin{pmatrix} x \\ 1 \end{pmatrix} \geq 0 \quad \Rightarrow f(x) \geq 0$$

Suppose \textcircled{1} holds

$$\text{consider } (x) \begin{pmatrix} A & b \\ b^T & c \end{pmatrix} \begin{pmatrix} x \\ t \end{pmatrix} \geq 0 \quad \forall x \in \mathbb{R}^n$$

$$= \begin{pmatrix} x^T A + tb^T & x^T b + tc \end{pmatrix} \begin{pmatrix} x \\ t \end{pmatrix} \geq 0$$

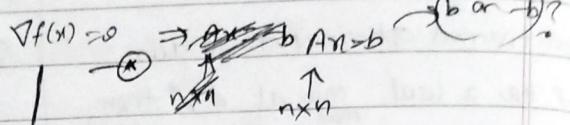
$$= x^T A x + tb^T x + x^T b t + t^2 c \geq 0$$

$$= x^T A x + t b^T x + t^2 b^T b + t^2 c \geq 0$$

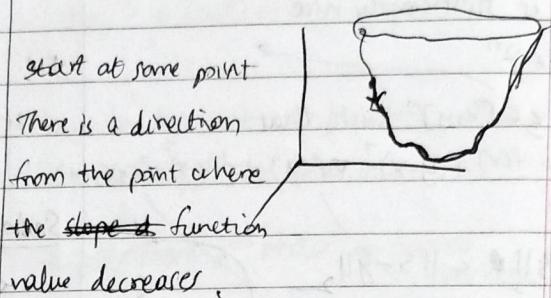
$$= t^2 \left(\left(\frac{x^T}{t} \right) A \left(\frac{x}{t} \right) + 2 b^T \left(\frac{x}{t} \right) + c \right) \geq 0$$

$$= t^2 f\left(\frac{x}{t}\right) \geq 0$$

$$\min_{x \in \mathbb{R}^n} f(x)$$



Analytical method is harder to use
 Descent method: To solve \textcircled{1} iteratively



find the descent

direction d_k

find the step

size t_k

- Stopping criterion

$$\| \nabla f(x_k) \| \leq \epsilon$$

$$x_{k+1} = x_k + t_k d_k$$

$$\begin{cases} f(x_{k+1}) < f(x_k) \\ f(x_{k+1}) \leq f(x_k) \end{cases} \quad = f(x_k + t_k d_k)$$

$\nabla f(x_k) \neq 0$

2 things to find
 1. t_k
 2. d_k

Lemma: (Sufficient decrease lemma)

Suppose $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is sufficiently nice and $0 \neq d \in \mathbb{R}^n$ is a descent direction, and $\alpha \in (0, 1)$

Then $\exists \epsilon > 0$ st.

$$f(x + \alpha d) \leq f(x)$$

$$0 \geq f(x) - f(x + \alpha d) \geq -\alpha [\nabla f(x), d]$$

$\forall \alpha \in [0, \epsilon]$

Sufficient decrease lemma

Aug 08

d_k is a descent direction

$$f'(x_k + d_k) = \nabla f(x_k) \cdot d_k \stackrel{\text{def}}{< 0}$$

stopping criterion

$$\|\nabla f(x_k)\| \leq \epsilon$$

$$\|x_{k+1} - x_k\| \leq \epsilon$$

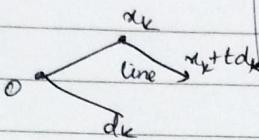
on finding t_k and d_k :

Finding t_k is referred to as line search

$$\text{line } g(t) = f(x_k + t d_k)$$

ways of finding t_k

• constant step size



$$t_k = \epsilon \quad \forall k$$

Pros Cons

• Simplicity

• At some k , might be good,
at other k , might be bad.

• Exact line search

$$\min_{t \geq 0} g(t) \quad (\text{now this is itself is an optimization problem})$$

We want $g(t) = f(x_k + t d_k)$ to be a "good" function

- Not be non-convex
- Not be ugly looking function that makes everything more complex.

$$t_k = \arg \min_{t \geq 0} g(t)$$

• Backtracking

$$\alpha, \beta \in (0, 1) \quad S > 0$$

start with $t_k = S$

while $f(x_k) - f(x_k + t_k d_k) \leq -t_k \alpha \|\nabla f(x_k)\| d_k$

$$t_k \rightarrow \beta \cdot t_k \quad (\text{we want sufficient decrease lemma to hold})$$

In other words $t_k = \beta^{i_0}$'s where i_0 is the smallest integer +ve integer such that

$$f(x_k) - f(x_k + t_k d_k) \geq -t_k \alpha \|\nabla f(x_k)\| d_k$$

holds

adv

Simple

Theoretically ideal

Practical

disadv

Sometimes not useful

Not always practical

constant step size
exact line search
Backtracking

Example:

$$f(x) = x^T A x \quad \text{where } A \text{ is P.D. and symmetric}$$

Find expression for t_k

$$g(t) = f(x_k + t d_k)$$

$$f(x) = x^T A x$$

$$g(t) = (x_k + t d_k)^T A (x_k + t d_k)$$

$$= (x_k^T + t d_k^T) A (x_k + t d_k)$$

$$= x_k^T A x_k + t^2 d_k^T A d_k$$

$$+ t x_k^T A d_k + t d_k^T A x_k$$

$$g'(t) = \cancel{2 x_k^T A d_k} + 2 t d_k^T A d_k + x_k^T A d_k + d_k^T A x_k$$

$$g'(t) = 0 \Rightarrow t_k = \frac{-(x_k^T A d_k + d_k^T A x_k)}{2 d_k^T A d_k}$$

$$= -\frac{x_k^T A d_k}{d_k^T A d_k}$$

Finding $d_k \rightarrow f'(x_k, d)$

$$\min_{d \in \mathbb{R}^n} \nabla f(x_k) \cdot d$$

$$\|d\|_2 = 1$$

Consider

$$\nabla f(x_k) \cdot d \geq -\|\nabla f(x_k)\|_2 \|d\|_2 \quad (\text{Cauchy-Schwarz})$$

$$\nabla f(x_k) \cdot d \geq -\|\nabla f(x_k)\| \quad \forall 0 \neq d \in \mathbb{R}^n$$

lower bound where $\|d\|_2 = 1$

$$\text{Set } d = -\frac{\nabla f(x_k)}{\|\nabla f(x_k)\|_2}$$

$$d_k^T A x_k = (-1) \underbrace{\| \nabla f(x_k) \|_2}_{} \quad x_k^T A d_k = -1 \underbrace{\| \nabla f(x_k) \|_2}_{} \\ = (-1) \underbrace{\| \nabla f(x_k) \|_2}_{} = \underbrace{\| \nabla f(x_k) \|_2}_{} = \boxed{1}$$

$$\text{Set } dz = \frac{\nabla f(x_k)}{\|\nabla f(x_k)\|_2}$$

$$f\left(\frac{x_{k+1} - \nabla f(x_k)}{\|\nabla f(x_k)\|_2}\right) = \nabla f(x_k) \cdot \left(\frac{-\nabla f(x_k)}{\|\nabla f(x_k)\|}\right)$$

$$= -\|\nabla f(x_k)\| \quad \text{(*)}$$

$\cap c_i$ is also convex
 $i \in I$

$\lambda_1 c_1 + \lambda_2 c_2 + \dots + \lambda_n c_n$ is also convex,

$c_1 \times c_2 \times c_3 \times \dots \times c_n \times \dots$ is also convex,
where $A \in \mathbb{R}^{m \times n}$

$A(c) = \{Ax \mid x \in C\}$ is also convex
↑
convex.

From (1), (**) optimal direction is

$$d_k = \frac{-\nabla f(x_k)}{\|\nabla f(x_k)\|}$$

Gradient (or steepest descent) method:

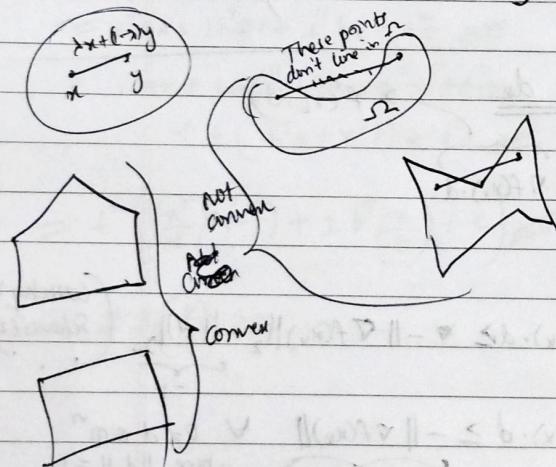
$$x_{k+1} = x_k - t_k \nabla f(x_k)$$

t_k can be obtained via line search

Convex sets

$S \subseteq \mathbb{R}^n$ is convex

If $[\lambda x + (1-\lambda)y] \in S \forall (x,y) \in S$ and $\forall \lambda \in [0,1]$



Results: Suppose $\{c_i\}_{i \in I}$ is a set of convex sets where $c_i \subseteq \mathbb{R}^2$ and I is finite or infinite.

$$(x_1, x_2) \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix} \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix}$$

$$\begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix} \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix} = \begin{pmatrix} a_{11} + a_{21} & a_{12} + a_{22} \\ a_{11} & a_{22} \end{pmatrix} \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix}$$

$$\sum \beta_i (a_{11} + a_{21}) = \sum \beta_i a_{11}$$