

Section 7. Linear Regression

1. Linear Regression
2. Least Squares
3. Matrix Calculus

Instructor: He Wang

Department of Mathematics

Northeastern University

➤ **House Price Example:**

Consider the 59 single family residential houses sold in Newton, MA in Dec. 2020.
(Data downloaded from www.redfin.com)

house 1:

BEDS	BATHS	LOCATION	SQUARE_FEET	LOT_SIZE	YEAR_BUILT	PRICE
3	3	Newton	2969	15014	1967	1090000
3	2.5	Newton	1566	5582	1922	805000
4	2.5	Newton Corner	2532	6273	1953	905000
7	4.5	Newton Center	6748	26607	1902	2660000
4	4	West Newton	4200	20446	2007	1925000
4	2.5	Newton	2232	3966	1870	965000
2	1.5	Newton Corner	1344	5559	1851	775000
3	2.5	Newton	2898	12420	1943	1250000
2	2	West Newton	1729	4171	1953	815000
6	3	West Newton	3149	12616	1953	900000
5	3.5	West Newton	4000	12006	1912	1800000
4	3.5	West Newton	6430	30600	1920	3550000
4	1.5	Auburndale	1750	8222	1893	885000
2	2	Newton	840	5548	1955	630000
...

New house : x_1 x_2 x_3

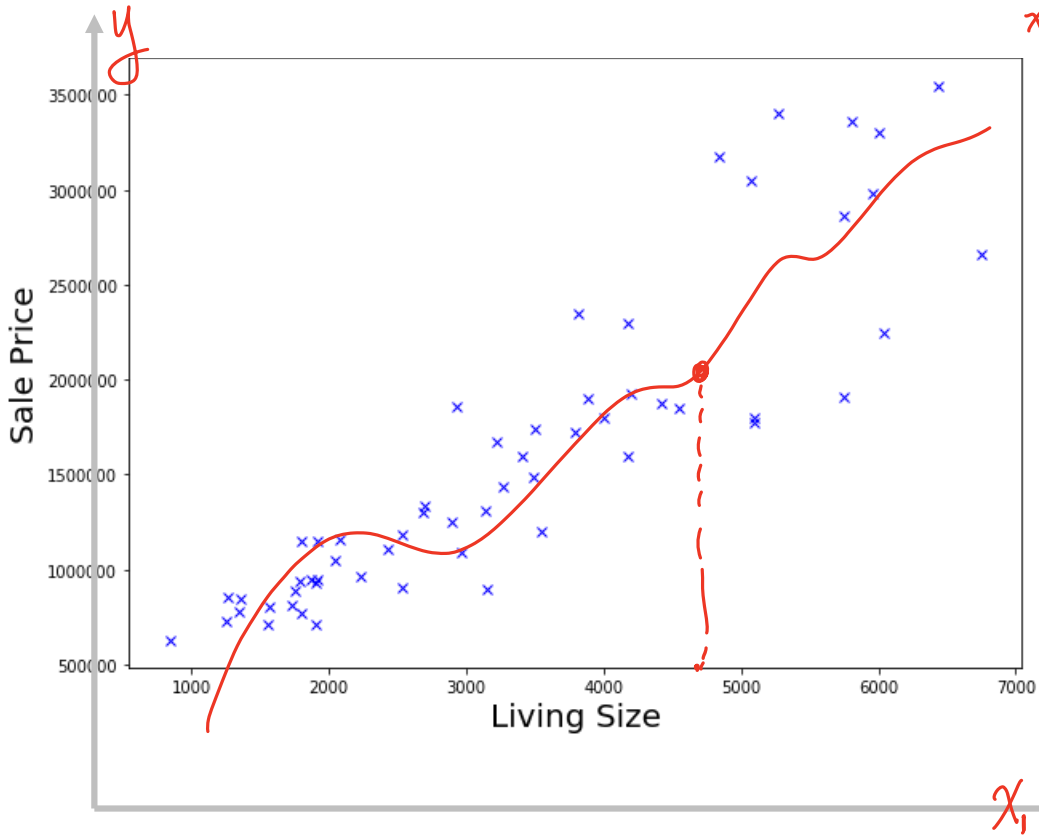
γ

predict

➤ Predict house price via living size (square_feet)

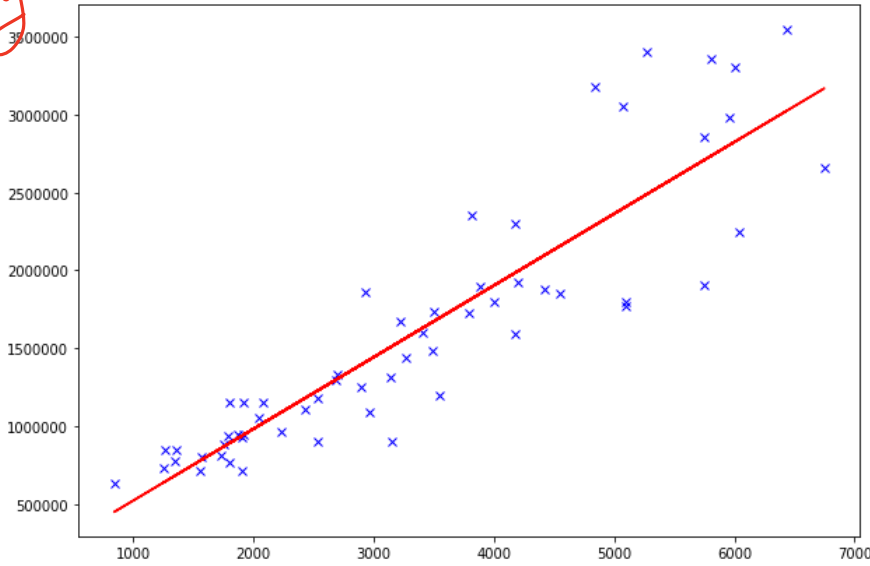
Input: a dataset that contains n samples. $(\vec{x}^{(i)}, y^{(i)})$, $i = 1, \dots, n$

Task: if a house has x_1 square feet, predict its price?



x_1	y
SQUARE_FEET	PRICE
2969	1090000 $y^{(1)}$
1566	805000
2532	905000
6748	2660000
4200	1925000
2232	965000
1344	775000
2898	1250000
1729	815000
3149	900000
4000	1800000
6430	3550000
1750	885000
840	630000
...	...

y



affine / "linear"

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1$$

e.g

$$h(x) = 3 + 4x_1$$

x_1

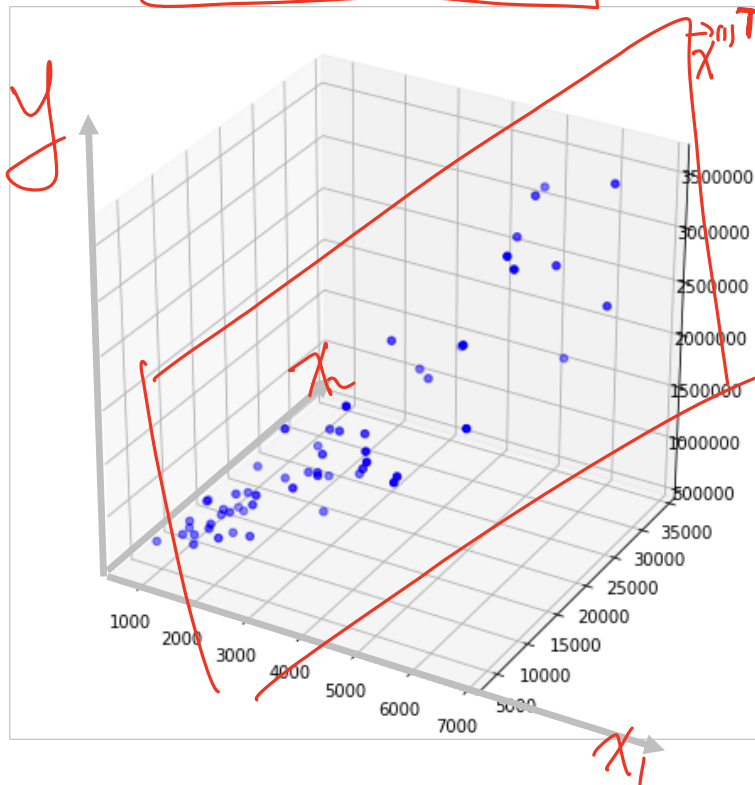
➤ Predict house price.

Input: a dataset that contains n samples $(\vec{x}^{(i)}, y^{(i)})$, $i = 1, \dots, n$

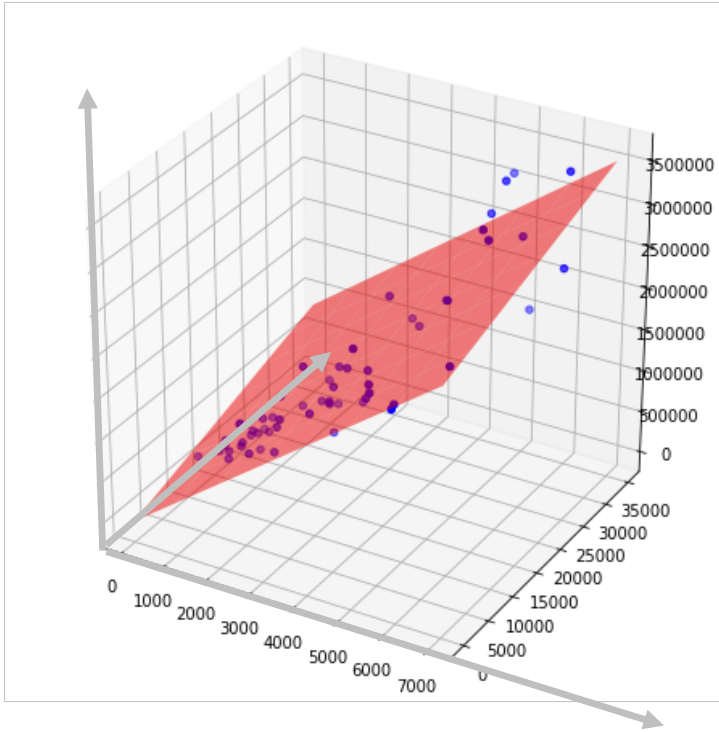
$$\vec{x}^{(1)} = \begin{bmatrix} 2969 \\ 15014 \end{bmatrix}$$

Task: if a house has x_1 (ft²) living size and x_2 (ft²) lot size, predict its price?

Model: $h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$



	x_1	x_2	y
	SQUARE_FEET	LOT_SIZE	PRICE
house 1	2969	15014	1090000
,	1566	5582	805000
,	2532	6273	905000
,	6748	26607	2660000
,	4200	20446	1925000
	2232	3966	965000
	1344	5559	775000
	2898	12420	1250000
	1729	4171	815000
	3149	12616	900000
	4000	12006	1800000
	6430	30600	3550000
	1750	8222	885000
	840	5548	630000
...



$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$$

➤ **Linear Regression (Parametric Method)**

Input: a dataset that contains n samples

$$D = \{(\vec{x}^{(i)}, y^{(i)}), \quad i = 1, \dots, n\}$$

$$\vec{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ \vdots \\ \vdots \\ x_d \end{bmatrix} \in \mathbb{R}^d \xrightarrow{h} y \in \mathbb{R}^1$$

“Model”
Assumption: linear model

$$\begin{aligned} h_{\theta}(\vec{x}) &= \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_d x_d \\ &= \sum_{i=0}^d \theta_i x_i = \vec{\theta} \cdot \vec{x} \\ &= \vec{\theta}^T \vec{x} = \vec{x}^T \vec{\theta} \end{aligned}$$

	x_1	x_2	...	x_d	y
SQUARE_FEE
T	LOT_SIZE	BEDS	BATHS	PRICE	
2969	15014	3	3	1090000	
1566	5582	3	2.5	805000	
2532	6273	4	2.5	905000	
6748	26607	7	4.5	2660000	
4200	20446	4	4	1925000	
2232	3966	4	2.5	965000	
1344	5559	2	1.5	775000	
2898	12420	3	2.5	1250000	
1729	4171	2	2	815000	
3149	12616	6	3	900000	

$$\vec{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix} \quad \vec{x} = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_d \end{bmatrix} \in \mathbb{R}^{d+1}$$

Data give us

$$h_{\theta}(\vec{x}^{(i)}) = y^{(i)} \text{ for } i = 1, \dots, n$$

$$h_{\theta}(\vec{x}^{(1)}) = y^{(1)}$$

Data and linear assumption implies

$$(\vec{x}^{(i)})^T \vec{\theta} = y^{(i)} \text{ for } i = 1, \dots, n$$

Matrix Notation:

$$\mathbf{X} \vec{\theta} = \vec{y}$$

$$\begin{bmatrix} \vec{x}^{(1)T} \\ \vdots \\ \vec{x}^{(n)T} \end{bmatrix} \vec{\theta} = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix}$$

Data Matrix:

Target vector:

Parameter vector:

$$\mathbf{X} = \begin{bmatrix} 1 & x_1^{(1)} & x_2^{(1)} & \dots & x_d^{(1)} \\ 1 & x_1^{(2)} & x_2^{(2)} & \dots & x_d^{(2)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(n)} & x_2^{(n)} & \dots & x_d^{(n)} \end{bmatrix} = \begin{bmatrix} \vec{x}^{(1)T} \\ \vec{x}^{(2)T} \\ \vdots \\ \vec{x}^{(n)T} \end{bmatrix}$$

$$\vec{y} = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix}$$

$$\vec{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix}$$

Data $(\vec{x}^{(i)}, y^{(i)})$ (X, \vec{y})
 $i=1, \dots, n$

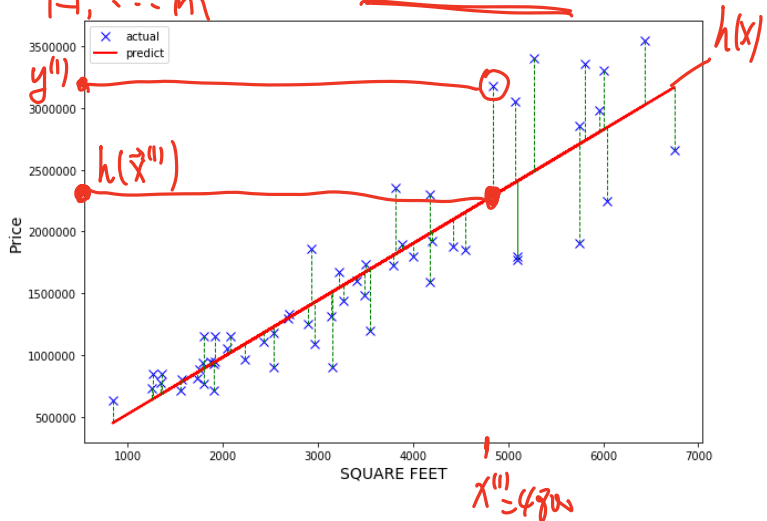
➤ Evaluate the model:

$h(x)$

Prediction Vector:

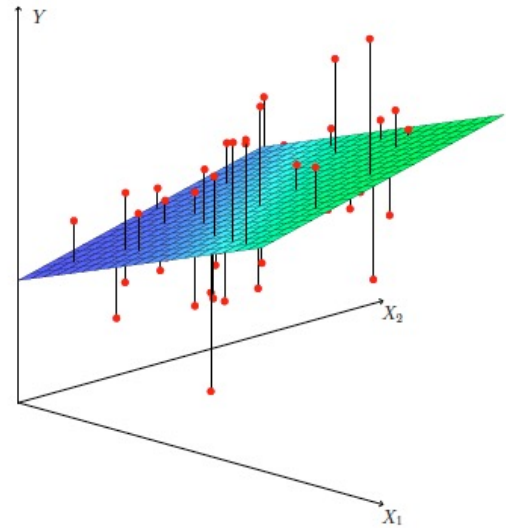
$$h(\mathbf{X}) = \begin{bmatrix} h(\vec{x}^{(1)}) \\ h(\vec{x}^{(2)}) \\ \vdots \\ h(\vec{x}^{(n)}) \end{bmatrix} - \vec{y} = \begin{bmatrix} -y^{(1)} \\ \vdots \\ -y^{(n)} \end{bmatrix}$$

true



Difference Vector

$$h(\mathbf{X}) - \vec{y} = \begin{bmatrix} h(\vec{x}^{(1)}) - y^{(1)} \\ h(\vec{x}^{(2)}) - y^{(2)} \\ \vdots \\ h(\vec{x}^{(n)}) - y^{(n)} \end{bmatrix}$$



➤ Cost/Loss Functions

- **Mean Absolute Error**

$$L(\vec{\theta}) = \frac{1}{n} \sum_{i=1}^n |h_{\theta}(\vec{x}^{(i)}) - y^{(i)}| = \frac{1}{n} \|h(\mathbf{X}) - \vec{y}\|_1$$

l₁-norm

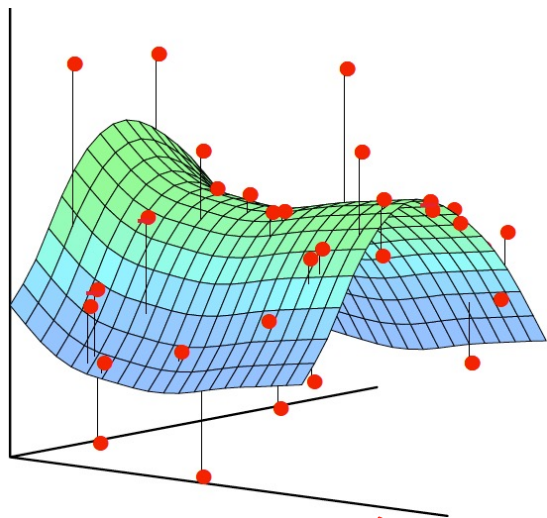
- **Mean Residual Sum of Squares**

$$L(\vec{\theta}) = \frac{1}{n} \sum_{i=1}^n (h_{\theta}(\vec{x}^{(i)}) - y^{(i)})^2 = \frac{1}{n} \|h(\mathbf{X}) - \vec{y}\|_2^2$$

- **Residual Sum of Squares (RSS):**

$$RSS(\vec{\theta}) := \sum_{i=1}^n (h_{\theta}(\vec{x}^{(i)}) - y^{(i)})^2 = \|h(\mathbf{X}) - \vec{y}\|_2^2$$

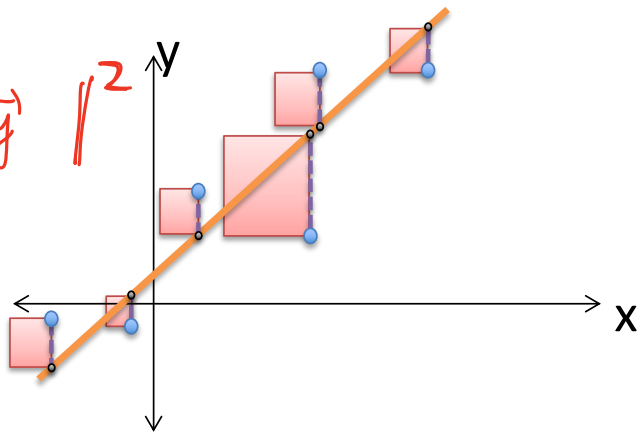
$$= \|h_{\theta}(\mathbf{X}) - \vec{y}\|_2^2$$



$$\vec{u} \cdot \vec{v} = u_1 v_1 + \dots + u_n v_n$$

$$\|\vec{u}\|^2 = u_1^2 + \dots + u_n^2$$

Picture Interpretation of RSS



➤ **Review: Inner Product, Norm, Metric on vector spaces V**

Let V be a real vector space. For example, V is a subspace of \mathbb{R}^n .

Definition (Inner Product). An inner product on V is a binary function

$$\langle -, - \rangle: V \times V \rightarrow \mathbb{R}$$

such that for vectors $\vec{u}, \vec{v}, \vec{w} \in V$ and a scalar $c \in \mathbb{R}$, the following hold:

- (1.) $\langle \vec{u}, \vec{v} \rangle = \langle \vec{v}, \vec{u} \rangle$
- (2.) $\langle \vec{u} + \vec{v}, \vec{w} \rangle = \langle \vec{u}, \vec{w} \rangle + \langle \vec{v}, \vec{w} \rangle$
- (3.) $\langle c\vec{u}, \vec{v} \rangle = c\langle \vec{v}, \vec{u} \rangle$
- (4.) $\langle \vec{u}, \vec{u} \rangle \geq 0$
- (5.) $\langle \vec{u}, \vec{u} \rangle = 0$ if and only if $\vec{u} = \vec{0}$

We call V an **inner product space** with inner product $\langle -, - \rangle$.

Example: Dot product on \mathbb{R}^n .

Example: Weighted dot product on \mathbb{R}^n .

$$\langle \vec{u}, \vec{v} \rangle_W := \vec{u}^T W \vec{v}$$

Here, W is a positive-definite symmetric matrix

Definition (Norm). Let V be a real vector space. A norm on V is a function

$$\|-\|: V \rightarrow \mathbb{R}$$

such that for vectors $\vec{u}, \vec{v} \in V$ and a scalar $c \in \mathbb{R}$, the following hold:

(1.) $\|\vec{u}\| \geq 0$

(2.) $\|\vec{u}\| = 0$ if and only if $\vec{u} = \vec{0}$

(3.) $\|c\vec{u}\| = |c| \|\vec{u}\|$

(4.) The triangle inequality $\|\vec{u} + \vec{v}\| \leq \|\vec{u}\| + \|\vec{v}\|$

We call V a normed space with norm $\|-\|$.

Example: l_2 -norm induced by dot product.

Example: l_p -norm on \mathbb{R}^n

$$\|\vec{v}\|_p = \sqrt[p]{|v_1|^p + \dots + |v_n|^p}$$

Example: l_p -cost function $L(\vec{\theta}) = \underbrace{\|h_{\theta}(\mathbf{X}) - \vec{y}\|_p}_p^p$

Definition (Metric). Let S be a set. A metric(distance) on S is a binary function

$$d: S \times S \rightarrow \mathbb{R}$$

such that for vectors $\vec{u}, \vec{v}, \vec{w} \in S$ and a scalar $c \in \mathbb{R}$, the following hold:

- (1.) $d(\vec{u}, \vec{v}) = d(\vec{v}, \vec{u})$
- (2.) $d(\vec{u}, \vec{v}) = 0$ if and only if $\vec{u} = \vec{v}$
- (3.) $d(\vec{u}, \vec{w}) \leq d(\vec{u}, \vec{v}) + d(\vec{v}, \vec{w})$

We call S a metric space metric function d .

Examples:

1. If S is a vector space, metric is equivalent to norm.

2. The **discrete metric on S** , where $d(x, y) = 0$ if $x = y$ and

$d(x, y) = 1$ otherwise.

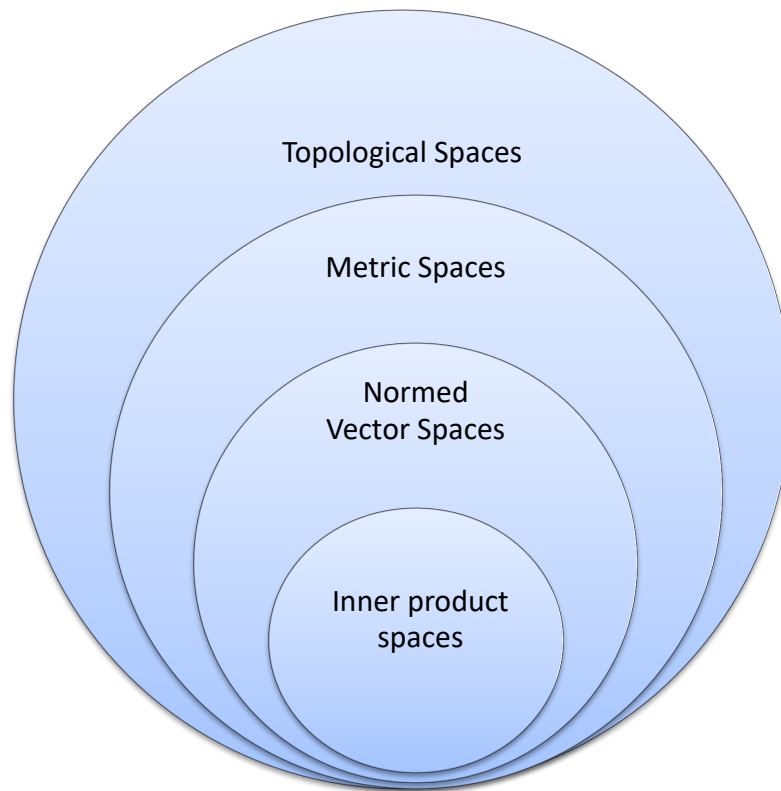
3. The **positive real numbers** with distance function $d(x, y) = |\log(y/x)|$ is a metric space.

Given a distance function d on the label space \mathcal{C}^n

- Cost/Loss Function: $L(\vec{\theta}) = d(h_{\vec{\theta}}(\mathbf{X}), \vec{y})$

prediction true value

End of review.



Ex: $h_{\theta}(\vec{x}) = \theta_0 + \theta_1 x_1 + \dots + \theta_d x_d$

➤ Minimize Cost/Loss Functions

$$(X, \vec{y})$$

$$X = \begin{bmatrix} \vec{x}^{(1)T} \\ \vdots \\ \vec{x}^{(n)T} \end{bmatrix} \quad \vec{y} = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix}$$

- Given labeled Data $\mathcal{D} = \{(\vec{x}^{(i)}, y^{(i)}), i = 1, \dots, n\}$
- Assumption(Model): $h_{\theta}(-)$
- Cost/Loss Function: $L(\vec{\theta}) = d(h_{\theta}(X), \vec{y})$ for some distance d on the label space.
∴ error!
or $\|h_{\theta}(X) - \vec{y}\|$ or $\langle h_{\theta}(X) - \vec{y}, h_{\theta}(X) - \vec{y} \rangle$

Goal: Find $\vec{\theta}$ to minimize the cost $L(\vec{\theta})$

Equivalently, find $\underset{\vec{\theta}}{\operatorname{argmin}} L(\vec{\theta}) = \{\vec{\theta} \text{ such that } L(\vec{\theta}) \text{ is minimized}\}$

$n \times d$

➤ Minimize Cost/Loss Functions (linear regression)

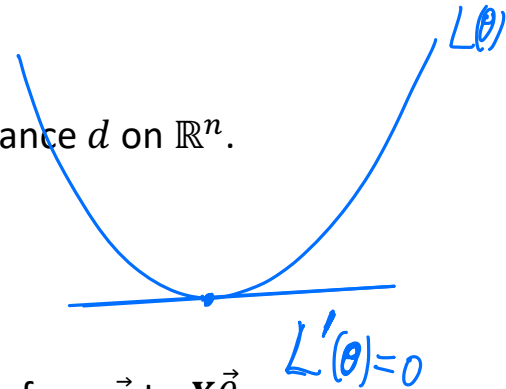
(X, \vec{y})

• Given labeled **Data** $\mathcal{D} = \{(\vec{x}^{(i)}, y^{(i)}), i = 1, \dots, n\}$, where $y^{(i)} \in \mathbb{R}$

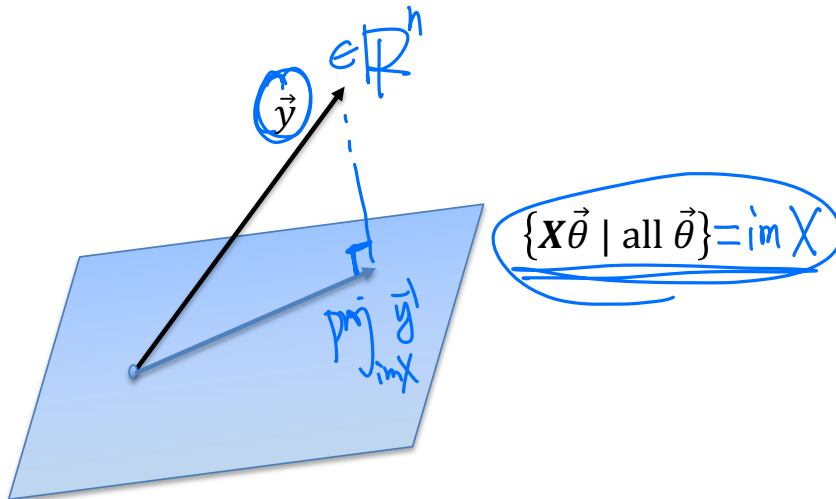
• Assumption(**Linear Model**): $h_{\theta}(\vec{x}) = \vec{x}^T \vec{\theta}$

• Cost/Loss Function: $L(\vec{\theta}) = d(h_{\theta}(\mathbf{X}), \vec{y})$ for some distance d on \mathbb{R}^n .

• Find argmin $L(\vec{\theta}) = \{\vec{\theta} \text{ such that } L(\vec{\theta}) \text{ is minimized}\}$



Minimize the cost is the same as minimize the distance from \vec{y} to $\mathbf{X}\vec{\theta}$



general

- If we the norm/distance is induced by an **inner product**, then the solution of $\underset{\vec{\theta}}{\operatorname{argmin}} L(\vec{\theta})$ is



$$\mathbf{X} \vec{\theta} = \operatorname{Proj}_{\operatorname{im}(\mathbf{X})} \vec{y}$$

example

- If we the norm is induced by **dot product**, the cost function is the residual sum of squares:

$$L(\vec{\theta}) = \operatorname{RSS}(\vec{\theta}) = \|h_{\theta}(\mathbf{X}) - \vec{y}\|^2 = \|\mathbf{X}\vec{\theta} - \vec{y}\|^2 = \sum_{i=1}^n \left((\vec{x}^{(i)})^T \vec{\theta} - y^{(i)} \right)^2$$

then the solution of $\underset{\vec{\theta}}{\operatorname{argmin}} L(\vec{\theta})$ is obtained by solving

$$\mathbf{X}^T \mathbf{X} \vec{\theta} = \mathbf{X}^T \vec{y}$$

This is called the **normal equation** of $\mathbf{X}\vec{\theta} = \vec{y}$

Lemma: Rank $X^T X = \text{Rank } X$

- If rank $X = \underline{d + 1}$, then $X^T X$ is invertible.

\times
 $n \times \underline{(d+1)}$

In this case, the solution for the normal equation is

arg min $\vec{\theta}$ = $\vec{\theta}^{\text{opt.}} = (X^T X)^{-1} X^T \vec{y}$

- If rank $X = d + 1$ and $X = QR$ where Q is an orthogonal matrix and R is an upper triangular matrix, then the solution for the normal equation is

$$\vec{\theta} = R^{-1} Q^T \vec{y}$$

Remarks:

$$\text{if } W=I \rightarrow \vec{u}^T \vec{v}$$

- If we the norm $\|\vec{u}\|_W := \langle \vec{u}, \vec{u} \rangle_W$ is induced by **weighted inner product** $\langle \vec{u}, \vec{v} \rangle_W := \vec{u}^T W \vec{v}$, the cost function is the:

$$L(\vec{\theta}) = \|\mathbf{X}\vec{\theta} - \vec{y}\|_W^2$$

The solution of $\underset{\vec{\theta}}{\operatorname{argmin}} L(\vec{\theta})$ is by solving the weighted normal equation

$$\mathbf{X}^T \mathbf{W} \mathbf{X} \vec{\theta} = \mathbf{X}^T \mathbf{W} \vec{y}$$

$$\text{if } W=I \text{ then } \mathbf{X}^T \mathbf{X} \vec{\theta} = \mathbf{X}^T \vec{y}$$

$$\text{if } W = \begin{pmatrix} c_1 & & \\ & \ddots & \\ & & c_n \end{pmatrix} \quad c_i > 0$$

- If the norm is not induced by inner product (e.g., the l_p -norm), then finding $\underset{\vec{\theta}}{\operatorname{argmin}} L(\vec{\theta})$ is a hard optimization question. We will need to use gradient descent or Newton's method to minimize the cost $L(\vec{\theta})$.

➤ **Matrix Calculus**

Ex: Assume $h_{\vec{\theta}} = \vec{x}^T \vec{\theta} = \vec{\theta}^T \vec{x}$ $h_{\vec{\theta}}(X) = X \vec{\theta}$

COST:

$$RSS(\vec{\theta}) = \|\mathbf{X}\vec{\theta} - \vec{y}\|^2 = (\mathbf{X}\vec{\theta} - \vec{y})^T (\mathbf{X}\vec{\theta} - \vec{y})$$

$$= (\vec{\theta}^T \mathbf{X}^T - \vec{y}^T) (\mathbf{X}\vec{\theta} - \vec{y})$$

$$\begin{aligned} & \vec{y}^T (\mathbf{X}\vec{\theta}) \\ &= (\mathbf{X}\vec{\theta})^T \vec{y} \\ &= \vec{\theta}^T \mathbf{X}^T \vec{y} \end{aligned}$$

Ex:

$$\begin{bmatrix} \theta_0 & \theta_1 \end{bmatrix} \begin{bmatrix} 1 & 3 \\ 3 & 2 \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix}$$

$$= \theta_0^2 + 2\theta_1^2 + 6\theta_0\theta_1$$

$$= \vec{\theta}^T \mathbf{X}^T \mathbf{X} \vec{\theta} - \vec{\theta}^T \mathbf{X}^T \vec{y} - \vec{y}^T \mathbf{X} \vec{\theta} + \vec{y}^T \vec{y}$$

$$= \underbrace{(\vec{\theta}^T \mathbf{X}^T \mathbf{X} \vec{\theta})}_{\text{Quadratic}} - \underbrace{2\vec{y}^T \mathbf{X} \vec{\theta}}_{\text{Linear}} + \underbrace{\vec{y}^T \vec{y}}_{\text{Constant}}$$

1x1 nx nx (n)

Data = (X, \vec{y})

Constant

$$\begin{bmatrix} 2 & 3 \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix} = 2\theta_0 + 3\theta_1$$

To minimize $RSS(\vec{\theta})$, we need to find critical points.

$$\begin{bmatrix} \] \] \end{bmatrix} = \in \mathbb{R}$$

$$J(\theta) = RSS(\vec{\theta}) : \mathbb{R}^{d+1} \rightarrow \mathbb{R}$$

$$\vec{\theta} \rightarrow RSS(\vec{\theta})$$

$$\left. \begin{aligned} \frac{\partial J}{\partial \theta_0} &= 0 \\ \frac{\partial J}{\partial \theta_1} &= 0 \end{aligned} \right\} \Rightarrow 2\theta_0 + 6\theta_1 + 2 = 0$$

$$\frac{\partial}{\partial \theta_i} = 0 \Rightarrow 4\theta_1 + 6\theta_0 + 3 = 0$$

➤ Matrix Calculus Math

(x_1, x_2, \dots, x_n)

Definitions. (Gradient/Partial derivative)

(1) If $f: \mathbb{R}^n \rightarrow \mathbb{R}$, the gradient of f is defined to be

$$\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} = \vec{x}$$

derivative,

$$\frac{\partial f}{\partial \vec{x}} := \nabla_{\vec{x}} f := \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix}$$

$$f(x_1, x_2) = f\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right)$$

$$= x_1^2 + 6x_1x_2 + 4x_2^2$$

$$\nabla f = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \end{bmatrix} = \begin{bmatrix} 2x_1 + 6x_2 \\ 6x_1 + 8x_2 \end{bmatrix}$$

(2) If $f: \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$, the gradient of f is defined to be

$$\begin{bmatrix} x_{11} & \dots & x_{1n} \\ \vdots & & \vdots \\ x_{m1} & \dots & x_{mn} \end{bmatrix} = X \rightarrow f(X)$$

$$\frac{\partial f}{\partial X} := \nabla_X f := \begin{bmatrix} \frac{\partial f}{\partial x_{11}} & \dots & \frac{\partial f}{\partial x_{1n}} \\ \vdots & \ddots & \vdots \\ \frac{\partial f}{\partial x_{m1}} & \dots & \frac{\partial f}{\partial x_{mn}} \end{bmatrix}$$

$$F: \mathbb{R}^2 \rightarrow \mathbb{R}^3$$

(3) If $F: \mathbb{R}^n \rightarrow \mathbb{R}^m$, the **derivative** of F is defined to be

$$\begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} = \vec{x} \rightarrow F(\vec{x}) = \begin{bmatrix} f_1(\vec{x}) \\ \vdots \\ f_m(\vec{x}) \end{bmatrix} \quad \frac{\partial F}{\partial \vec{x}} := \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_m}{\partial x_1} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_1}{\partial x_n} & \dots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

Ex: $F(\vec{x}) = \begin{bmatrix} x_1^2 + x_2^2 \\ x_1 x_2 \\ 3x_1 + 4x_2 \end{bmatrix} = \begin{bmatrix} f_1 \\ f_2 \\ f_3 \end{bmatrix} \in \mathbb{R}^3$

$$\frac{\partial F}{\partial \vec{x}} = \begin{bmatrix} 2x_1 & 2x_2 & 3 \\ x_2 & x_1 & 4 \end{bmatrix}$$

$n \times m$
 2×3

Remark: The above notation is called **denominator layout** notation, which is a generalization of the gradient.

There is another **numerator layout** convention for the derivative of F , which is the transpose of the denominator layout

"transpose"

$$\begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \dots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

Theorem(Linear Property)

Let f and g be functions and c be a real number.

$$(1) \nabla(f + g) = \nabla f + \nabla g$$

$$(2) \nabla(cf) = c\nabla f$$

Here ∇ can be $\nabla_{\vec{x}}$, or ∇_x , or $\frac{\partial}{\partial \vec{x}}$

$$f: \mathbb{R}^n \rightarrow \mathbb{R}^1$$

• **Proposition:** If $f(\vec{x}) = \vec{b}^T \vec{x}$, then $\nabla f = \vec{b}$.

$$\parallel \\ b_1 x_1 + \dots + b_n x_n$$

$$\nabla f = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix} = \begin{bmatrix} b_1 \\ \vdots \\ b_n \end{bmatrix} = \vec{b}$$

$$F: \mathbb{R}^n \rightarrow \mathbb{R}^n \quad \text{if } A \text{ } \underline{m \times n}$$

Proposition: If $F(\vec{x}) = A\vec{x}$, then $\frac{\partial F}{\partial \vec{x}} = A^T$

$(n \times n)$

$$A \vec{x}' = \begin{pmatrix} a_{11} & \dots & a_{1m} \\ \vdots & & \vdots \\ a_{m1} & \dots & a_{mn} \end{pmatrix} \begin{pmatrix} x_1 \\ \vdots \\ x_m \end{pmatrix}$$

Proof:

1. Write down $A\vec{x}$ explicitly using a_{ij} and x_j
2. Write down $\frac{\partial F}{\partial \vec{x}}$ explicitly as partial derivatives as definition.

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

Proposition: (quadratic function):

If $f(\vec{x}) = \vec{x}^T A \vec{x}$, then $\nabla f = (A^T + A)\vec{x}$

e.g. $A = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$

$$f(\vec{x}) = \begin{bmatrix} x_1 & x_2 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$$

$$= x_1^2 + 5x_1x_2 + 4x_2^2$$

~~Proof:~~ $f(\vec{x}) = \vec{x}^T A \vec{x} = \sum_{i=1}^n \sum_{j=1}^n a_{ij} x_i x_j$

$$\nabla f = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \vdots \\ \frac{\partial f}{\partial x_n} \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^n a_{i1} x_i + \sum_{j=1}^n a_{1j} x_j \\ \vdots \\ \sum_{i=1}^n a_{in} x_i + \sum_{j=1}^n a_{nj} x_j \end{bmatrix} = A^T \vec{x} + A \vec{x}$$

A

$(n \times n)$

Ex: If $A^T = A$, then $\nabla f = 2A\vec{x}$

Example: $J(\vec{\theta}) = \text{RSS}(\vec{\theta}) = \vec{\theta}^T \mathbf{X}^T \mathbf{X} \vec{\theta} - 2\vec{y}^T \mathbf{X} \vec{\theta} + \vec{y}^T \vec{y}$

\mathbf{X}
 $n \times (d+1)$

$$\begin{bmatrix} \frac{\partial J}{\partial \theta_0} \\ \frac{\partial J}{\partial \theta_1} \\ \vdots \\ \frac{\partial J}{\partial \theta_d} \end{bmatrix} = \nabla_{\vec{\theta}} J = \underline{2\mathbf{X}^T \vec{\theta}} - 2\mathbf{X}^T \vec{y} = \mathbf{0}$$

$$\mathbf{X}^T \mathbf{X} \vec{\theta} = \mathbf{X}^T \vec{y}$$

normal equation!

critical points

$$\begin{bmatrix} 2x_1x_2 \\ 3x_1x_3 \\ x_1x_3 \end{bmatrix}$$

$$f \cdot s = e^x \cdot \sin x$$

Theorem: (Product Rule) (denominator layout)

Suppose $F: \mathbb{R}^n \rightarrow \mathbb{R}^m$ and $G: \mathbb{R}^n \rightarrow \mathbb{R}^m$ and $H = G^T F$. Then,

$$\begin{matrix} \begin{matrix} \text{=} \\ \text{=} \end{matrix} \begin{bmatrix} F_1 \\ \vdots \\ F_m \end{bmatrix} & \begin{bmatrix} G_1 \\ \vdots \\ G_m \end{bmatrix} & \begin{matrix} \text{=} \\ \text{=} \end{matrix} \end{matrix} \quad \frac{\partial H}{\partial \vec{z}} = \frac{\partial G}{\partial \vec{z}} F + \frac{\partial F}{\partial \vec{z}} G$$

$n \times 1 \quad n \times m \quad m \times 1$

$F^T G: \mathbb{R}^n \rightarrow \mathbb{R}$

Pf: $H = G^T F = [g_1 \dots g_m] \begin{bmatrix} f_1 \\ \vdots \\ f_m \end{bmatrix} = g_1 f_1 + \dots + g_m f_m$

$$\begin{aligned} \frac{\partial H}{\partial z_i} &= \left(\frac{\partial g_1}{\partial z_i} f_1 + \frac{\partial f_1}{\partial z_i} g_1 \right) + \dots + \left(\frac{\partial g_m}{\partial z_i} f_m + \frac{\partial f_m}{\partial z_i} g_m \right) \\ &= \frac{\partial G}{\partial z_i} F + \frac{\partial F}{\partial z_i} G \quad \text{for each } i=1,2, \dots \end{aligned}$$

Chain Rule. Assume that $\vec{Y} \in \mathbb{R}^n$ is a vector depending on $\vec{X} \in \mathbb{R}^m$, and \vec{X} depends on some $\vec{Z} \in \mathbb{R}^q$. Then

$$\frac{\partial \vec{Y}}{\partial \vec{Z}} = \frac{\partial \vec{X}}{\partial \vec{Z}} \frac{\partial \vec{Y}}{\partial \vec{X}} \quad \frac{df}{dt} = \frac{df}{du} \cdot \frac{du}{dt}$$

second derivative

Definition: The Hessian matrix of $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is defined by

$$H(f) := \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix}$$

e.g.

$$f(\vec{x}) = x_1^2 + 3x_1x_n$$

Proposition: If $f(\vec{x}) = \vec{b}^T \vec{x}$, then $H(f) = 0$.

$$\underline{\underline{=}} \quad \underline{\underline{=}} \quad \nabla f = 2A\vec{x}$$

Proposition: If $f(\vec{x}) = \vec{x}^T A \vec{x}$, for a symmetric matrix A , then $H(f) = 2A$

$$\underline{\underline{=}} \quad \underline{\underline{=}} \quad A = A^T$$

Example: $J(\vec{\theta}) = \text{RSS}(\vec{\theta}) = \vec{\theta}^T \mathbf{X}^T \mathbf{X} \vec{\theta} - 2\vec{y}^T \mathbf{X} \vec{\theta} + \vec{y}^T \vec{y}$

$$H(J) = 2\mathbf{X}^T \mathbf{X}$$

symmetric, positive-semidefinite

Theorem: (Second Derivative Test)

If $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is smooth, then a critical point $\vec{a} \in \mathbb{R}^n$ (i.e., $\nabla f(\vec{a}) = \vec{0}$) is

- (1) a local minimum if $H(f(\vec{a}))$ is positive definite; \rightarrow all eigenvalues of $H(f(\vec{a})) > 0$
- (2) a local maximum if $H(f(\vec{a}))$ is negative definite; $\leftarrow < 0$
- (3) a saddle point if $H(f(\vec{a}))$ contains positive and negative eigenvalues;
- (4) there is no conclusion for the other cases. zero eigenvalues

Examples:

$$f(x, y) = x^2 + y^2$$

$$H(f) = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix} \text{ positive definite}$$

$$f(x, y) = -x^2 - y^2$$

$$f(x, y) = x^2 - y^2$$

$$H = \begin{bmatrix} 2 & \\ & -2 \end{bmatrix}$$

$$\nabla f = \begin{bmatrix} 2x \\ 2y \end{bmatrix} = 0 \Rightarrow \text{critical point} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$



$$f(x, y) = -x^4$$

$$f(x, y) = x^4$$

$$f(x, y) = x^2 + y^3$$

Example: $f(x, y) = (x + y)(xy + xy^2)$

<https://www.geogebra.org/3d/bjsj7erx>



Definition: A function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is **convex** if

$$f(\lambda \vec{u} + (1 - \lambda)\vec{v}) \leq \lambda f(\vec{u}) + (1 - \lambda)f(\vec{v}) \quad \text{for any } 0 \leq \lambda \leq 1$$

Definition: A set C is **convex** if and only if

$$\vec{u}, \vec{v} \in C \Rightarrow \lambda \vec{u} + (1 - \lambda)\vec{v} \in C \quad \text{for any } 0 \leq \lambda \leq 1$$

Theorem: A function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is smooth and convex if and only if

$$f(\vec{u}) - f(\vec{v}) \geq \nabla f \cdot (\vec{u} - \vec{v}) \quad \text{for all } \vec{u}, \vec{v}$$

- Any local minimum of a convex function is also a global minimum.
- If its Hessian $H(f(\vec{x}))$ is everywhere positive semi-definite, then f is convex.