# Multiple Regression

## Given

$$\mathcal{D} = \{x_i, y_i\}_{i=1}^{n}$$

$x_i \in \mathbb{R}^d$        Multiple independant variable

$y_i \in \mathbb{R}$        Single dependant variable

$$y_i = f(x_i)$$

## Objective:

$$\hat{f}: \mathbb{R}^d \to \mathbb{R}$$

## Hypothesis

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \ldots + \beta_d x_d + \varepsilon$$

$\beta_i \Rightarrow$ slope of respective axis     Random error

$$\hat{y} = w_0 + w_1 x_1 + w_2 x_2 + \ldots + w_d x_d + e$$

error in the model

(data points)
i = 1 to n
j = 1 to d
(feature)

$$Y = X w + e$$

$$w = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_d \end{bmatrix}_{(d+1) \times 1} \qquad Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}_{1 \times n}$$

$$X = \begin{bmatrix} 1 & x_{11} & x_{12} & x_{13} & \ldots & x_{1d} \\ 1 & x_{21} & x_{22} & x_{23} & \ldots \ldots & x_{2d} \\ \vdots & & & & \\ 1 & x_{n1} & x_{n2} & x_{n3} & \ldots & x_{nd} \end{bmatrix}_{n \times (d+1)}$$

## Find w

$$SSE = \min \sum_{i=1}^{n} (y_i - \hat{y}_i)^2$$

$$\min SSE = \sum_{i=1}^{n} (y_i - (w_0 + w_1 x_{i1} + w_2 x_{i2} + \ldots + w_d x_{id}))^2$$

UCOP

$$\frac{\partial L}{\partial w_0} = \sum_{i=1}^{n} 2\left(y_i - \left(w_0 + w_1 x_{i1} + w_2 x_{i2} + \ldots w_d x_{id}\right)\right)(-1) = 0$$

$$\sum_{i=1}^{n} y_i = \sum_{i=1}^{n} \left(w_0 + w_1 x_{i1} + w_2 x_{i2} + \ldots w_d x_{id}\right)$$

$$\sum y_i = n w_0 + w_1 \sum_{i=1}^{n} x_{i1} + w_2 \sum_{i=1}^{n} x_{i2} + \ldots w_d \sum x_{id}$$

$$\frac{\partial L}{\partial w_1} = \sum_{i=1}^{n} 2\left(y_i - \left(w_0 + w_1 x_{i1} + w_2 x_{i2} + \ldots w_d x_{id}\right)\right)(-x_{i1}) = 0$$

$$\sum y_i x_{i1} = w_0 \sum x_{i1} + w_1 \sum x_{i1}^2 + w_2 \sum x_{i1} x_{i2} + \ldots w_d \sum x_{i1} x_{id}$$

$$\frac{\partial L}{\partial w_2} = \sum_{i=1}^{n} 2\left(y_i - \left(w_0 + w_1 x_{i1} + w_2 x_{i2} + \ldots w_d x_{id}\right)\right)(-x_{i2}) = 0$$

$$\sum y_i x_{i2} = w_0 \sum x_{i2} + w_1 \sum x_{i1} x_{i2} + w_2 \sum x_{i2}^2 + \ldots w_d \sum x_{id}^2$$

$$\frac{\partial L}{\partial w_d} = \sum_{i=1}^{n} 2\left(y_i - \left(w_0 + w_1 x_{i1} + w_2 x_{i2} + \ldots + w_d x_{id}\right)(-x_{id}\right) = 0$$

$$\sum y_i x_{id} = w_0 \sum x_{id} + w_1 \sum x_{i1} x_{id} + \ldots \ldots w_d * \sum x_{id}^2$$

$$X^T = \begin{bmatrix} 1 & 1 & 1 & \ldots \ldots & 1 \\ x_{11} & x_{21} & \ldots \ldots & x_{n1} \\ x_{12} & x_{22} & \ldots \ldots & x_{n2} \\ \vdots & \vdots & & \vdots \\ x_{1d} & x_{2d} & \ldots \ldots & x_{nd} \end{bmatrix}_{(d+1) \times n}$$

$$X^T X = \begin{bmatrix} n & \sum_{i=1}^{n} x_{i1} & \sum_{i=1}^{n} x_{i2} & \ldots & \sum_{i=1}^{n} x_{id} \\ \sum_{i=1}^{n} x_{i1} & \sum_{i=1}^{n} x_{i1}^2 & \sum_{i=1}^{n} x_{i1} x_{i2} & \ldots & \sum_{i=1}^{n} x_{i1} x_{id} \\ \sum x_{i1}^2 & \sum x_{i1} x_{i2} & \sum x_{i2}^2 & \ldots & \sum x_{i2} x_{id} \\ \vdots & & & \\ \sum x_{id} & \sum x_{id} x_{i1} & \sum x_{id} x_{i2} & & \sum x_{id}^2 \end{bmatrix}$$

$$x^Ty = \begin{bmatrix} \sum y_i \\ \sum x_{i1} y_i \\ \sum x_{i2} y_i \\ \vdots \\ \sum x_{id} y_i \end{bmatrix}_{(d+1) \times 1}$$

$$x^Txw = x^Ty$$
$$w = (x^Tx)^{-1}x^Ty$$

$$w = (x^Tx)^{-1} x^Ty$$

| Cost | 7 | 3 | 3 | 4 | 6 | 7 |
|------|----|----|----|----|----|----|
| Distance | 560 | 220 | 340 | 80 | 150 | 330 |
| Time | 16.68 | 11.50 | 12.03 | 14.88 | 13.75 | 18.11 |

$$x = \begin{bmatrix} 1 & 7 & 560 \\ 1 & 3 & 220 \\ 1 & 3 & 340 \\ 1 & 4 & 80 \\ 1 & 6 & 150 \\ 1 & 7 & 330 \end{bmatrix} \qquad x^T = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 & 1 \\ 7 & 3 & 3 & 4 & 6 & 7 \\ 560 & 220 & 340 & 80 & 150 & 330 \end{bmatrix}$$

$$y = \begin{bmatrix} 16.68 \\ 11.50 \\ 12.03 \\ 14.88 \\ 13.75 \\ 18.11 \end{bmatrix}$$

$$x^Tx = \begin{bmatrix} 6 & 30 & 1680 \\ 30 & 1680 & 9130 \\ 1680 & 9130 & 615400 \end{bmatrix} \qquad (x^Tx)^{-1} = \begin{bmatrix} 1.60723 & -0.25 & -6\times10^{-4} \\ -0.25 & 0.0698 & -3\times10^{-6} \\ -6\times10^{-4} & -3\times10^{-6} & 8.6\times10^{-6} \end{bmatrix}$$

$$x^Ty = \begin{bmatrix} 86.95 \\ 456.14 \\ 25190.2 \end{bmatrix} \qquad w = (x^Tx)^{-1}x^Ty$$

$$w = \begin{bmatrix} 8.5655 \\ 1.1965 \\ -2.\times10^{-4} \end{bmatrix}$$

$$\hat{y} = 8.5655 + 1.1965 \, x_1 + (-2\times10^{-4}) \, x_2 + e.$$

Assumptions (No correlative variables)

* No multicolinearity

Measures

SSE
SST
SSR
$R^2$

$$\text{Adjusted } R^2 = 1 - \frac{(1-R^2)(n-1)}{(n-d-1)}$$

Increase d
Adjusted $R^2$ ↓se

$$= 1 - \frac{SSE/n-(d+1)}{SST/n-1}$$

i) Correlations

Check for each pair of features

Correlation    Covariance

$$-1 \le \sigma_{AB} = \frac{\Sigma_{AB}}{\sigma_A \sigma_B} \le +1$$

Standard Deviation

+1 = Positive correlated   (A ↑se then B ↑se)

0 = Independent.

Construct Correlation matrix

$$\begin{bmatrix} & \\ & \end{bmatrix}_{d \times d} \quad [d = \text{features}]$$

$$-0.7 \le \sigma_{ij} \le 0.7$$
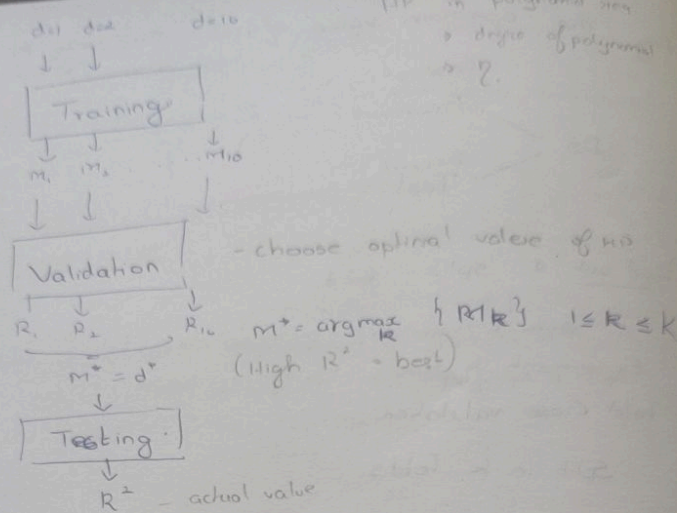
ii) Variance Inflation factor

$$VIF = \frac{1}{1-R^2}$$

$x_1, x_2, x_3, x_4 \ldots x_d \qquad Y$

$x_i x_j \Rightarrow$ keeping one as independent & other dependent
Construct linear reg model & find VIF

No of pairs = $dC_2$
(No of models to construct)

$$\begin{bmatrix} & \\ & \end{bmatrix}_{d \times d}$$

Highly correlation VIF is high ∴ remove 1 feature

HP in polynomial reg
  → degree of polynomial
  → ?

dil dil    dil0
 ↓  ↓       ↓

```
┌─────────────┐
│  Training   │
└─────────────┘
```
 ↓  ↓       ↓
$M_1$  $M_2$  ....$M_{10}$
 ↓  ↓       ↓

```
┌─────────────┐
│ Validation  │          - choose optimal value of HP
└─────────────┘
```
 ↓   ↓       ↓
$R_1$  $R_2$    $R_{10}$   $M^* = \underset{k}{argmax} \{ R[k] \}$   $1 \leq k \leq K$
  $M^* = d^*$        (High $R^2$ → best)
  ↓

```
┌─────────────┐
│  Testing    │
└─────────────┘
```
  ↓
$R^2$ — actual value

ERROR

- Training Error
- Test error / True error / Generalised error

    Traing Error $\begin{cases} \omega \\ min \; SSE \end{cases}$

    TE = 0 ; Fn passes through all data points
        TE = Large : Doesnot pass through many data point

    Test error ≥ 0 ; Test data points lie on Line.

  Bias = Training error
  Variance = Test error (How varied from Training error)
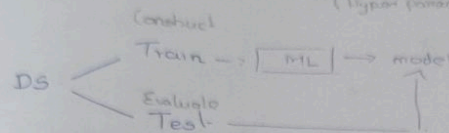                              (deviated)

  Bias                                                      Polynomial
                                                            degree high
- Low = model fits train data well / Complex model
- High = model doesnot fit train data well


  Variance

- Low = less deviation from Training error / Good model

| Var \ Bias | Low | High |
|---|---|---|
| Low | Ideal | X |
| High | Complex Over Fitting | Under fitting |

Finding the optimum value of H.P. → 3 way split
(Hyper parameter)

Construct

Train → | ML | → model

DS {

Evaluate
Test _____

Model Selection

• Hold out
• K fold Cross validation
• Leave - one - out
• 3 way split
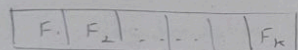
Hold out → split 2:1

Size (Train DS) > Size (Test)

Experiment k times & it' avg taken

### k fold Cross validation

Split in k folds

At i experiment i th is used as Test set other fold
used as Train set

Train set and Test set at any point must be disjoint

| F₁ | F₂ | . . . | . . . | Fₖ |

Every fold tested & remain used as Train set

Lastly find avg

### Leave one out (n Fold cross validation)
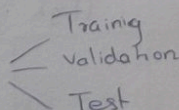
n folds
Each fold consist of 1 data point

Every point tested while other n-1 points kept as Trainning

n is large      -   Hold out
(Data set large)     k fold

Data set small   -   Leave one out

### 3 way split

Split in 3 data sets   { Training
                        Validation
                        Test

Reasons for overfit & Under fit

1) Insufficient Trainset
2) Noise / outlier

Increase data set
  - Augmentation — Do transformation
  - Smole analysis

Overfitting (overcome)

  * Ridge Regression / $L_2$ - Regularisation
  * Lasso      "      / $L_1$ -    "
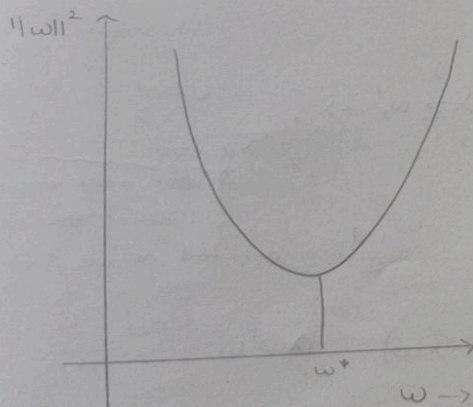  * Elastic net

→To all parametric mode

Find
$$\min_{w} \sum_{i=1} (y_i - \hat{y_i})^2 + \text{Regularization term}$$
$$+ \quad \|w\|^2$$

$$w^T = (w_0 \ldots w_d)$$

$$\|w\| = \sqrt{w_0^2 + w_1^2 + \cdots + w_d^2} \quad - L_2 \text{ Norm}$$



Ridge regularisation

Find w
$$\min \quad SSE + \lambda \cdot L_2 \text{ norm of } w$$

$$\Downarrow$$

Find w
$$\min \sum_{i=1} (y_i - \hat{y_i}) + \lambda \|w\|^2 \longrightarrow \text{square for easier derivation}$$
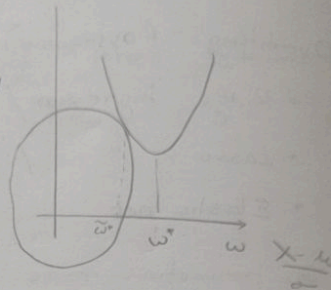
Find $w$

$$\min \sum_{i=1}^{n} (y^i - \hat{y_i})^2 + \lambda \sum_{i=1}^{d} w_j^2$$

Loss function : Continuous & Differentiable fn

Ridge Regularisation    D-dimensional sphere

Prerequisite (Not consider bins)

Data points to be normalised

Aand the center



Find $w$

$$\min \sum_{i=1}^{n} (y_i - \hat{y_i})^2$$

$\Downarrow$

Find $w$

$$\min \sum_{i=1}^{n} (y_i - w^T x_i)^2 \qquad \text{Convex fn}$$

} Constrained Convex optimisation problem

s.t.c

$$w_1^2 + w_2^2 + w_3^2 + \cdots + w_d^2 = 1 \qquad \text{Convex set}$$

$\Downarrow$

→ Lagrange multiplier → used to convert Constrained op to unconstrained op

Find $w$

$$\min \sum_{i=1}^{n} (y_i - w^T x_i)^2 + \lambda \left( \sum_{j=1}^{d} w_j^2 - 1 \right)$$

$\Downarrow$

$$\sum_{i=1}^{n} (y_i - w^T x)^2 + \lambda \sum_{j=1}^{d} w_j^2 - \lambda$$

U COP

→ $\lambda$ is const
I won't affect value of obj fn Hence removed.

$\boxed{Y_{n \times 1} \quad X_{n \times d} \quad W_{d \times 1} \quad W^T_{1 \times d}}$

$d+1$ to $d$ as we dont include 1

Find $w$   $\overbrace{\phantom{xxxxx}}$  $(Y - xw)^2$

$$L = \underset{n \times 1}{(Y - xw)^T} \; \underset{n \times d \; d \times 1}{(y - xw)} + \lambda w^T w$$

$\boxed{\begin{array}{c} \text{Matrix} \\ \text{Notation} \\ x^2 = x^T x \end{array}}$

$(y^T - w^T x^T)(y - xw)$  $_{1 \times n}$    $_{n \times 1}$

$$L = (y^T y) - \underline{y^T x w} - \underline{w^T x^T y} + w^T x^T x w + \lambda w^T w$$

Equal ( Transpose of one another)

$$= (x^T y) + 2 x^T w^T y + x^T w^T x + x w^T w$$

$$= y^T y - 2 w^T x^T y + w^T x^T x w + \lambda w^T w,$$

$$= y^T y - 2 w^T x^T y + w^2 x^T x + \lambda w^2$$

$$\frac{\partial L}{\partial w} = -2x^T y + 2wx^T x + 2\lambda w$$

$$= -x^T y + w(x^T x + \lambda I) w = 0$$

$$w(x^T x + \lambda I) w = x^T y$$

$$w = (x^T x + \lambda I)^{-1} x^T y$$

$\lambda = 0$    Similar to Multiple Regression.

Overfitting not Solved.

$w_j$ explains the importance of feature $x_{ij}$

Ridge regularisation $\Rightarrow$ solution to Multicolinearity

Ridge regularisation

$\therefore$ the importance of features $\downarrow$

## Lasso regularisation. (feature selection technique)

Min

SSE + L, Norm

$$\sum (y_i - \hat{y_i})^2 + \lambda \left( \sum_{j=1}^{d} (|w_j|) \right)$$

Obj = Convex fn

Constraint = Convex set.

$\therefore$ COP

$\Downarrow$

s.t.

Min $\sum_{j=1}^{n} (y_i - w^T x_i)^2$



s.t. $\sum_{j=1}^{d} |w_j| = 1$

$\Downarrow$

Min $\sum_{i=1}^{n} (y_i - w^T x_i)^2 + \lambda \left( \sum_{j=1}^{d} |w_j| - 1 \right)$

$\Downarrow$

Min $\sum (y_i - w^T x_i)^2 + \lambda \sum_{j=1}^{d} |w_j|$

$// w = 0$   feature not selected $\therefore$ used to select feature

$$w = (x^T x + \lambda I)^{-1} \cdot x^T y$$

Elastic Net

Both $L_1$ & $L_2$ norm added.

Classification (class variable - discrete)

Given

$$\mathcal{D} = \{ x_i, y_i \}_{i=1}^{n}$$

$$x_i \in R^d \; [\text{Discrete} \mid \text{Continuous}]$$

$$y_i \in \{ c_1, c_2 \}$$

$$f : x_i \to y_i$$

Objective

$$\hat{f} : R^d \to \{ c_1, c_2 \}$$

Logistic Regression - attempt to apply regression in classification

$$y_i = \text{discrete}$$

$$\mathcal{D} = \{$$

↓ LR (Min SSE
but we need line to divide
$w_0 + w_1 x_1 + \ldots w_0 x_d$ exactly)

↓

x   $y = c_1 / c_2$
LR   but gives continuous values
to $\Big\downarrow_{\text{change to discrete}}^{w^T x}$
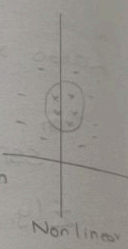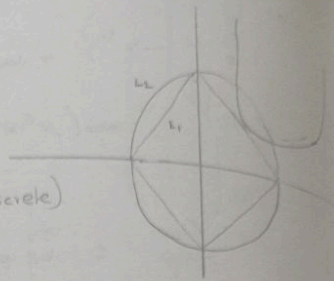
Sigmoid function

$$0 < \sigma(x) < 1 \qquad \text{for} \quad -\infty \leq x < +\infty$$

(Range of 0 to 1)

↓ threshold value 0.5

$$\hat{y} = \begin{cases} c_1 & \text{if } \sigma(w^T x) \geq 0.5 \\ c_2 & \text{o/w} \end{cases}$$
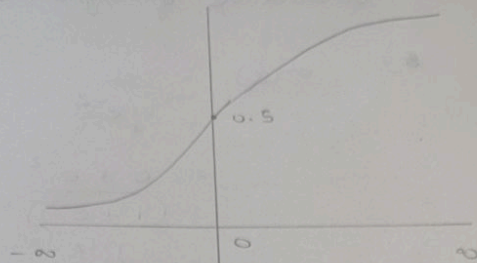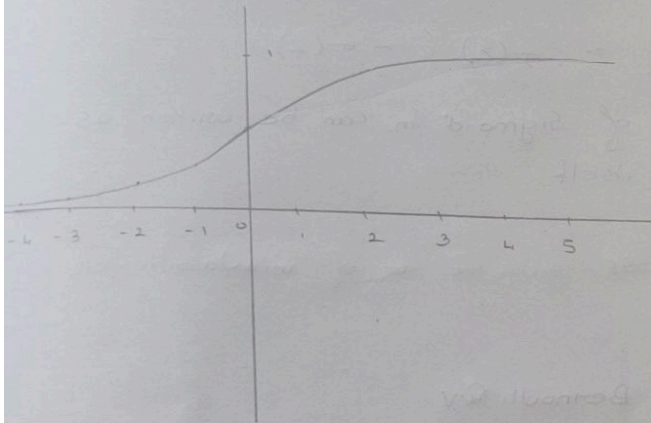
Sigmoid function                                      S shaped fn

$$\sigma(z) = \frac{1}{1 + e^{z}}$$

$\sigma(0) = 0.5$



Plot for $\sigma(x)$ ; $-5 \leq x \leq 5$



$\sigma(-1) = 0.26$
$\sigma(-2) = 0.11$
$\sigma(-3) = 0.04$
$\sigma(-4) = 0.01$
$\sigma(-5) = 0.006$
$\sigma(5) = 0.99$
$\sigma(4) = 0.98$
$\sigma(3) = 0.95$
$\sigma(2) = 0.88$
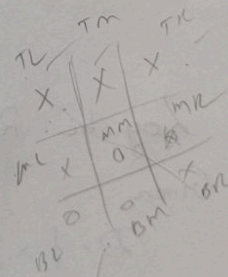$\sigma(1) = 0.73$

Properties of Sigmoid function

1) $0 < \sigma(z) < 1$     $\forall z \quad -\infty \leq z \leq \infty$

2) $\sigma(z) = \frac{1}{2}$ for $z = 0$

3) $z \to +\infty \quad \sigma(z)$ close to $1$

4) $z \to -\infty \quad \sigma(z) \quad '' \quad '' \quad 0$

5) if $\sigma(z) = \frac{1}{1+e^{-z}}$ , $1 - \sigma(z) = \frac{e^{-z}}{1+e^{-z}}$

## Derivative of Sigmoid Fn

**a)** $\sigma'(z) = \dfrac{1}{1+e^{-z}}$

$$= \frac{0 + e^{-z}}{(1+e^{-z})^2}$$

$$= \frac{e^{-z}}{(1+e^{-z})^2}$$

$$\sigma'(z) = \frac{1}{(1+e^{-z})} \cdot \frac{e^{-z}}{(1+e^{-z})}$$

$$= \sigma(z)(1 - \sigma(z))$$

derivative of Sigmoid fn can be written as Sigmoid fn itself

In classification SSE gives the no of misclassification (as it is Binary)

If $y$ is a Bernoulli R.V

$$f(y) = p^y (1-p)^{1-y} \qquad p = p(success)$$

$$y = 0 \qquad f(y) = 1 - p \text{ (failure)}$$

$$y = 1 \qquad f(y) = p \qquad \text{(success)}$$

## Logistic Regression

Given

$$D = \{x_i, y_i\}_{i=1}^{n}$$

$$x_i \in \mathbb{R}^d \quad (x_i \sim N(\mu, \sigma)) \text{ Gaussian}$$

Normal R.V (Continuous)

$$y_i \in \{c_1, c_2\}$$

$$f : x_i \rightarrow y_i$$

Objective

$$\hat{f} : \mathbb{R}^d \rightarrow \{c_1, c_2\}$$

Show that $\{0^n : n$ is prime$\}$ is not regular

Suppose $L$ is regular. By Pumping lemma

$\exists m$, $\forall \omega$ with $|\omega| \geq m$, $\exists \omega = xyz$ such that

$|xy| \leq m$, $|y| \geq 1$, $xy^k z \in L$ $\forall k \geq 0$

for any $m$,

choose, $\omega = 0^p$ where $p$ is prime & $p > m$.

Consider, $x = 0^i$

$\qquad y = 0^j$

$\qquad z = 0^{i-i-j}$, $xy^k z$.

$xy^k z$ where $\ell = |z| + |z| \in L$

$xy^2 z \in L$

$|xy^2 z|$ is prime

$|x| + |y^2| + |z|$ is prime

$|x| + |z| + 2|y|$ is prime.

$\ell + 2|y|$ is prime

$\ell(1 + |y|)$ is prime,

Note: $1 + |y| \neq 1$

By pumping lemma,

$xy^0 z \in L$

$|xy^0 z|$ is prime

$|x| + 0|y| + |z|$ is prime

$|x| + |z|$ is prime

$\ell$ is prime

$\ell \neq 1$

when $\ell \neq 1$ & $1 + |y| \neq 1$

then $\ell(1 + |y|)$ is not prime,

∴ Contradicts.

show that $\{0^n : n$ is not prime$\}$ not regular.

Suppose $L$ is regular, then $L^c$ is also regular.

$L^c : \{0^n : n$ is prime$\}$ is regular

a contradiction.

$L = \{a^n b^n : n \geq 0\}$ & $\cup$ $\{a^n b^m : n, m \geq 2\}$ is regular?

$L_1 \cup L_2$, if $L_1 \subseteq L_2$

$\Rightarrow L_2 \cup \{ab, \lambda\}$

- union of 2 regular language is regular.

$L = \{a^n b^n : n \neq 100\}$ is not regular.

$L_1 = \{a^n b^n : n \geq 0\}$

$L_2 = \{a^{100} b^{100}\}$ $(n = 100)$

$L_1 = L \cup L_2$

Since $L_1$ is regular (ASSUM)

$L_2$ is regular (finite)

contradict. $\therefore L_1 = L \cup L_2$ is also regular.

ST $L = \{\omega : |\omega|_a = |\omega|_b\}$ is not regular.

$L_1 = \{a^n b^n : n \geq 0\}$

$L_2 = a^* b^*$

$L_1 = L \cap L_2$

for CA1

→ formal language, operation on languages, Find $L^R, L^c, L^*,$

⊗→ pumping lemma

→ give a reg Ex, defn of DFA, NFA, reg Ex.

→ diff bet NFA & DFA

→ how to construct DFA?

→ product automata (starting with a & ending with bb)

→ equivalence, minimization of DFA.

⊗ conversion of NFA (with ∈, w/o ∈) to DFA

⊗ DFA to reg Ex, reg Ex to DFA.

→ Ardens lemma (no, state elimination / any one method)

→ closure properties. ($L_1$ is reg, how to prove $L_1 \cup L_2$ is reg, $L_1^c$ is reg, etc...)