Unsupervised learning → clustering
→ density estimation

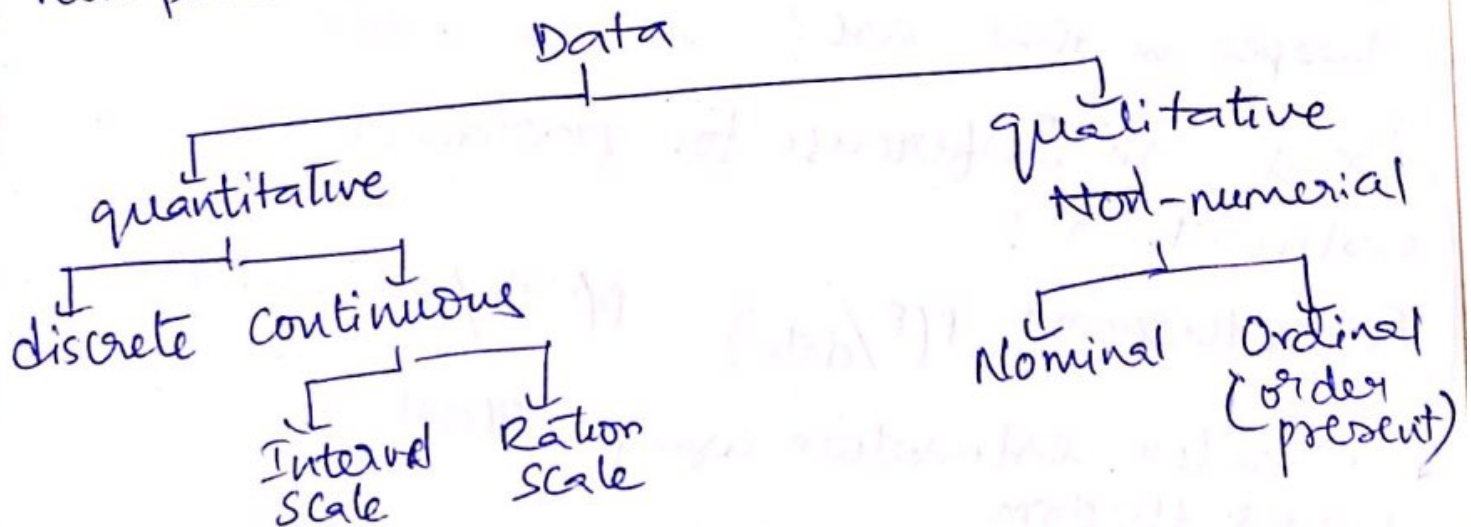Parametric        non-parametric

Parametric
ULearnip
① construct
gaussian mixture
models
② Use expectation
maximization algo

Clustering as a mixture of gaussians
"parametric (vs) non-parametric
distribution".

Different statistical distributions : (next page)

*If data is ordinal/interval based, only
non-parametric statistics can be used

Data
├── quantitative
│   ├── discrete
│   └── continuous
│       ├── Interval Scale
│       └── Ratio Scale
└── qualitative Non-numerical
    ├── Nominal
    └── Ordinal (order present)

Sequential learning
mapping input seq to output seq using
state machines. Hidden state seq present.

Active learning.

Theory of rational agency: (action selection theories)

*Density estimation ((how using deep generative NNs?))
→ Estimating probability density function of random variable in a population from sample's help.

q : Difference between probability density fn & probability distribution?
+n

q: what is maximum likelihood estimation?
→ finding the values of parameters that result in best fit curve.

* likelihood & loglikelihood

$$L(\mu, \sigma; data) = P(data; \mu, \sigma)$$

q when is least squares minimization same as max likelihood estimation? why does it & happen in that case? How

Bayesian Inference for parameter estimation:

Bayes theorem: $P(\theta/data) = \dfrac{P(data/\theta) \times P(\theta)}{P(data)}$

*Parameter estimation using Bayes theorem

"prior distribution".
$\theta \to$ set of parameters ($\theta = \{\mu, \sigma\}$ for gaussian distribution)
$P(\theta/data) \to$ posterior distribution
$P(\theta) \to$ prior distribution
$P(data) \to$ evidence & data = $\{y_1, y_2, \dots, y_n\}$
↳ normalizing const (helps making $\Sigma P(\theta/data) = 1$)

Can we use bayesian inference for classification problems? How? Is it used for discrete data / continuous or both?

Different statistics from the posterior distribution & their physical significance!
→ expected value ≡ mean
→ variance ⇒ uncertainity
→ mode ≠ MAP estimate.

"gaussian distribution is conjugate to itself wrt gaussian likelihood function."

(Latent Dirichlet Allocation algo) ✷

Markov Chain Monte Carlo methods → to Calculate posterior distribution

✷✷ Updating beliefs iteratively in real time. using bayesian inference → Kalman filter.
Prior acts as a regularizer here in bayesian inference.

q when does MAP estimate equals MLE?
⊕ Overfitting due to Bayesian priors ⟶ |pending|

Marginalization ⤵

$$P(x) = \int_y P(x, Y = y) \, dy.$$

What is discriminant analysis?
when is it used?
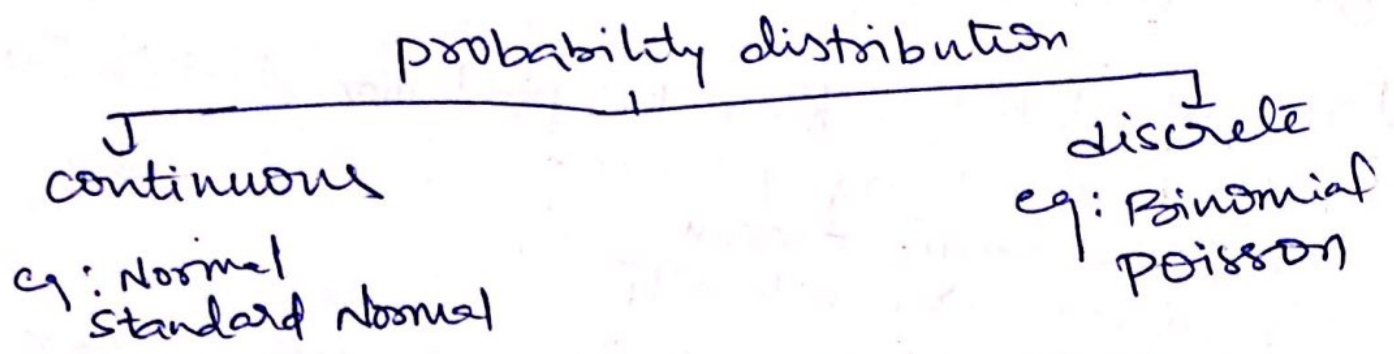when dependant variable is categorical & predictor / independant variable is interval.
develop discriminant fn as a linear combination of independant variables to descriminate between categories of dependant variable.

Discriminant analysis (vs) Analysis of variance.
(vs) regression analysis
Correlation is not causation. #

probability distribution

continuous                                    discrete
                                              eg: Binomial
eg: Normal                                       Poisson
   standard Normal

## MCMC methods

→ monte carlo simulations
→ markov chains (are memoryless)
# bell curve, law of large nos.
Markov → Non independant events may also
             conform to patterns
             (In long run, dist settle to pattern)
       * if events are subj to fixed prob.,
    interdependant events conform to average.

q. How can bayesian inference be used to
quantify uncertainity in predictions?

MCMC → Random sampli of parameters in
probabilistic space to approximate the posterior
distribution in ~~Bayesian~~ yesian inference.

~~where~~ can we use these posterior distributions?
How
→ quantifying uncertainity
→ comparing models
→ generati predictions

Central limit theorm & law of large nos
Covariance vs correlation vs causation

# Statistical distributions

q: data discrete/continuous?
q. symmetry of data & outliers scenario.
q. upper & lower limits of data
q. likelihood of occurence of extreme values

| Discrete distributions : | Continuous |
|---|---|
| → Binomial | → Normal |
| → Poisson | → Exponential |
| → Negative binomial | → logistic cauchy |
| → geometric | → gamma |
| → Bernoulli | → chisquared |

q. lets say we have a distribution which is generated by combining multiple commonly known distributions. How do we find those & seperate those distributions?

## Joint distributions:

## Discriminant analysis

quantifies uncertainty of estimated skill of model

quantifies uncertainty in single forecast

Exponential regression : $y = \alpha e^{\beta x}$ → Can be converted to linear

power regression : $y = \alpha x^{\beta}$

confidence interval vs prediction interval vs tolerance interval

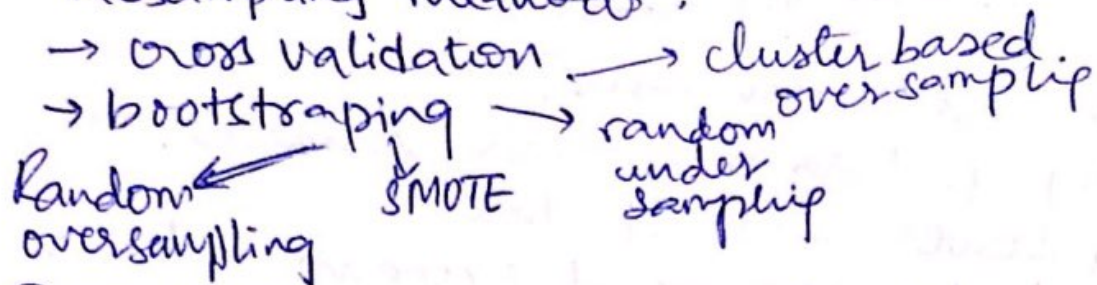q: Finding confidence interval ~~using~~ using different ~~means.~~ distribution

q. How to get confidence interval for mean given we have a sample?

## Methods to get prediction intervals :
→ bootstrap resampling
→ delta method
→ bayesian method
→ mean-variance estimation method

# Resampling methods:

→ cross validation ——→ cluster based
→ bootstraping ——→ random oversamplig
Random ←—— SMOTE    under samplig
oversamplling

Probability mass fn (vs) probability density fn

## Hypothesis testing:

→ null hypothesis
→ p value & critical value
→ traditional testy (vs) bayesian tesling

## Marginal distributions

X, Y are jointly distributed random variables

Entropy → list of generalized entropies
most widely used → Shannon's entropy

$$\phi_1(p) = -\sum_{i=1}^{n} P_i \log P_i$$

Decision tree → info gain → $\phi_{parent} - \sum w_i \times \phi_{child}$
         → gini index                    $w_i = \dfrac{n_{child}}{n_{parent}}$
                    → $1 - \sum_i P_i^2$

hinge loss → SVM

grid search → best combination of hyperparam-
     random search  eters for best fit

Skewness, kurtosis, coefficient of variation

Hyperparameter
   optimization

Data imputation methods : < → for cross-sectional datasets
                            → for time series data

Cases:
    → MCAR
    → MAR
    → NMAR

① Using mean / Median values :
  → poor results on encoded categorical features

② Using Mode
  → works with categorical features
  → can introduce bias

③ Using KNN algo.
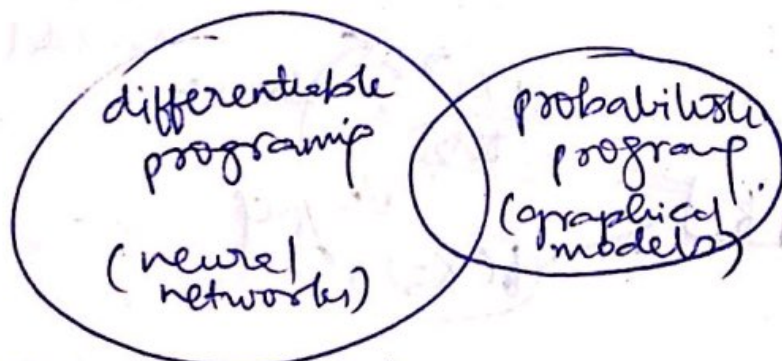  → impyute library
    KDTree
  KNN sensitive to outliers unlike SVM.

④ using multivariate imputation by chained eqn.
  → Multiple imputations

⑤ Using DL.

⑥ using stochastic regression/extrapolation & intrapolation

Uncertainity quantification

→ gaussian process (GP) models.
(multivariate problems)

differentiable programing (neural networks)

probabilistic program (graphical models)

gaussian process :

→ use prior knowledy to make predictions.
→ assign probability to diff ans possible for
fit a dataset, getting mean of prob dist to get
                                              most probable
→ incorporate confidence.                        ans.

- multivariate gaussian distribution :
  - → each random variable → normal distributed
  - → their joint dist → also normal
  - → $\mu$ → mean
    - $\Sigma$ → covariance matrix (symmetric & +ve semidefinite)
      - → gives $\sigma_i^2$ & $\sigma_{ij}$

Semi supervised learning → model trained on dataset → small portion labelled data
  ↳ majority unlabelled data .

- step 1 : cluster similar data into grps of similar data (unsup part)
- step2 : label remaing data in grp seeing the labelled data in the same grp.
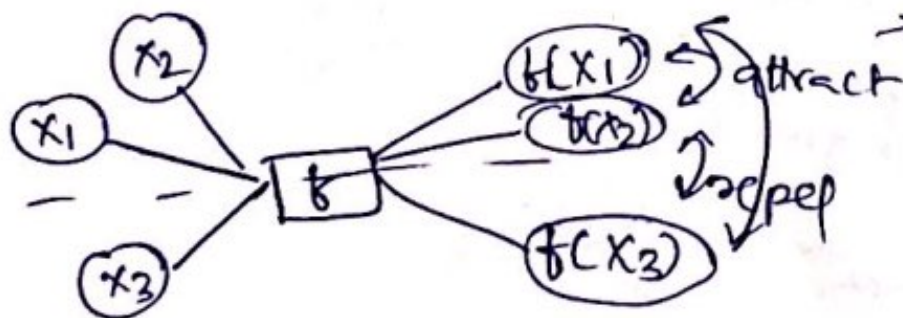
manifold assumption →

## contrastive learning :



adversial ML
  → data poisoning
  → evasion attack
  → model extractn

## Few short learning :
  ↳ meta learning problem

data level approach ← → parameter level approach ↳ limit parameter space

→ gans
→ data augmentation
→ using base datasets

N-way-k-shot-Classification :
learn how to learn to classify.
Meta-learning algo

## graph embedding methods & node embedding

→ learning multiple embeddings for a node

                    embedding
transductive ⟵ ─────────── → inductive

                    approaches
                                    → deep learning
factorization ⟵───────╲
methods            random
                    walks

★ → gaussian distribution based graph embedding
    (includes uncertainty estimation)

for :
→ node classification
→ link prediction
→ community detection

1) Matrix factorization:
   → using adjacency matrix → most simple method
   → using (locally linear embedding )(LLE):

   $$E_i = \sum_{j \in N_i} w_{ij} \times E_j$$

   $$\phi(E) = \sum_i \left( E_i - \sum_{j \in N_i} w_{ij} \times E_j \right)^2$$

2) HOPE :

   $$\phi(E) = \sum_{ij} \left( E_i \, E_j^T - S_{ij} \right)^2$$

   ⟵ similarity btw nodes i, j.
   (using Adamic/Adar
                ↑ similarity)

3) Deep walk

   $$A(x,y) = \sum_{u \in N(x) \cap N(y)} \frac{1}{\log |N(u)|}$$

word embedding techniques:

→ TF-IDF
⌐→ Word2vec ⟨ → Skip-gram
⌐→ GLOVE → Continuous bag of words
└→ BERT
→ bag of words

Gaussian process ⟵ (set of random variables)

$$x = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \sim N(\mu, \varepsilon) \qquad x \sim N(\mu_x, \varepsilon_{xx})$$

$$\varepsilon = Cov(x_i, x_j) = E[(x_i - \mu_i)(x_j - \mu_j)^T]$$

Marginalization & conditioning

$$P_{x,y} = \begin{bmatrix} x \\ y \end{bmatrix} \sim N(\mu, \varepsilon) = N\left( \begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix}, \begin{bmatrix} \varepsilon_{xx} & \varepsilon_{xy} \\ \varepsilon_{yx} & \varepsilon_{yy} \end{bmatrix} \right)$$

$$Y \sim N(\mu_y, \varepsilon_{yy})$$

$$P_x(x) = \int_y P_{x,y}(x,y) \, dy = \int_y P_{x/y}(x/y) \, P_y(y) \, dy$$

$$x/y \sim N\left( \mu_x + \varepsilon_{xy} \varepsilon_{yy}^{-1} (y - \mu_y), \ \varepsilon_{xx} - \varepsilon_{xy} \varepsilon_{yy}^{-1} \varepsilon_{yx} \right)$$

$$y/x \sim N\left( \mu_y + \varepsilon_{yx} \varepsilon_{xx}^{-1} (x - \mu_x), \ \varepsilon_{yy} - \varepsilon_{yx} \varepsilon_{xx}^{-1} \varepsilon_{xy} \right)$$

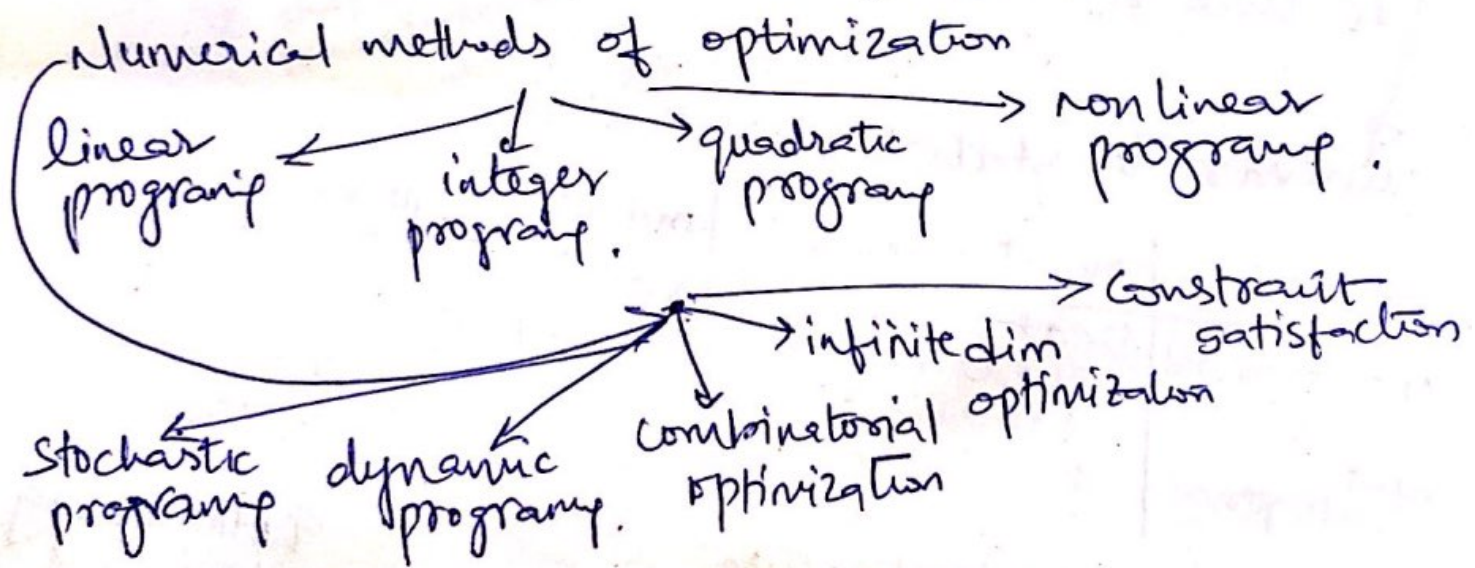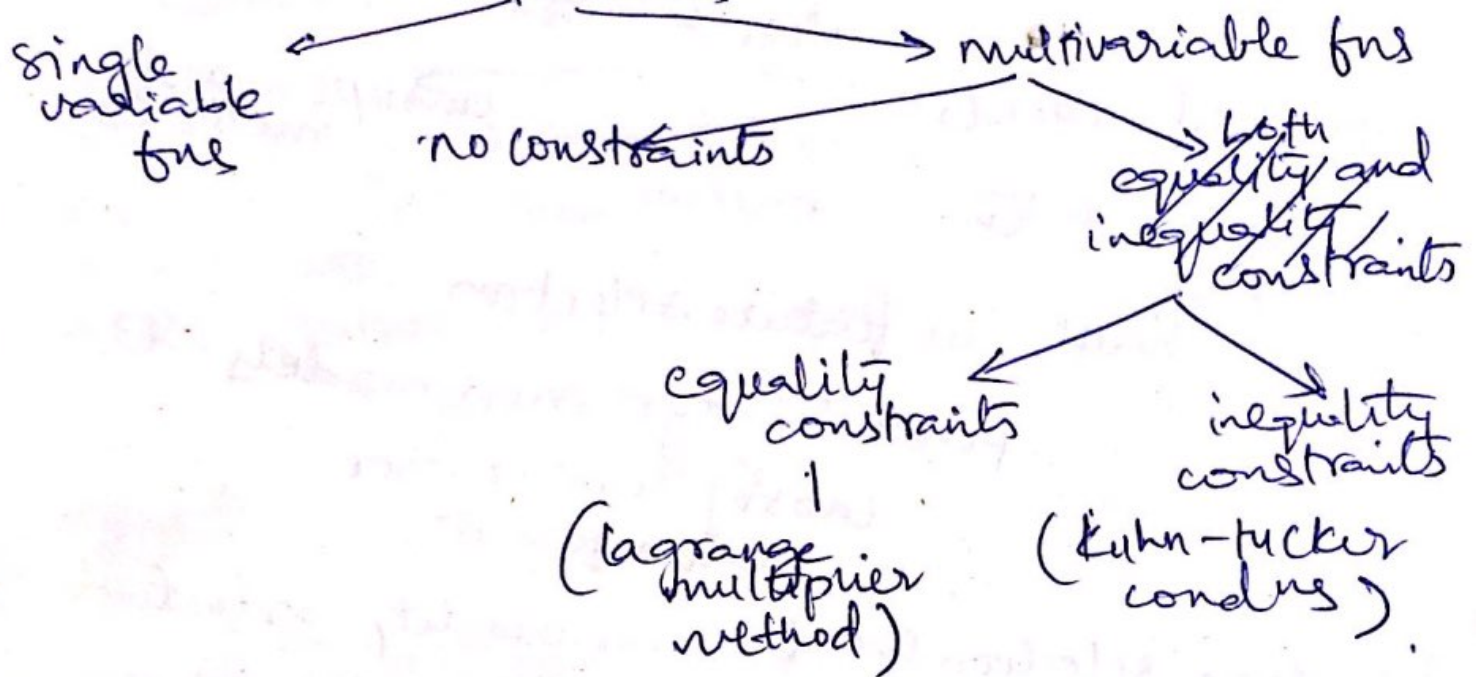each random variable has index $i$.
            ith dimension of ⤴
                    n dimensional multivariate
                                distribution.

# Optimization techniques

problems

single variable fns ← → multivariable fns

no constraints

both equality and inequality constraints

equality constraints ← → inequality constraints

|

(lagrange multiplier method)

(Kuhn-tucker condns)

Numerical methods of optimization

linear program ← integer program → quadratic program → non linear program

→ Constraint satisfaction

→ infinite dim optimization

combinatorial optimization

Stochastic program   dynamic program

advanced optimization techniques

Hill climbing ← genetic algorithms

(inheritance, mutation, selection, crossover/recombination)

# Correlation (diff from causation)

## Feature selection

Supervised methods ──────────── unsupervised methods

wrapper ──→ filter

### Built - in feature selection

penalized regression models

lasso ──→ decision tree
random forest

## Feature selection (vs) Dimensionality reduction

### univariate statistical measures

|  | output numerical | output categorical |
|---|---|---|
| input numerical | Pearson's | anova |
|  | Spearman's | kendall's |
|  | Anova | chi squared |
| input categorial | kendall's | mutual information |

### Categorical feature

q. Handling categorical data in ML

nominal ──→ ordinal

one hot encoding    target encoding    binary encoding    backward difference encoding    label encoding

## Handling missing data in time series :

→ Last observation carried forward
→ Next observation carried backward
→ Linear interpolation
→ Spline interpolation
→ when seasonality present, 1. deseasonalize
  2. interpolate 3. seasonalize

# Outlier analysis

univariate                                          multivariate

→ Normalize the data
→ Check z-score of each datapoint
→ if $x \notin [-3, 3]$, then outlier
→ using IQR
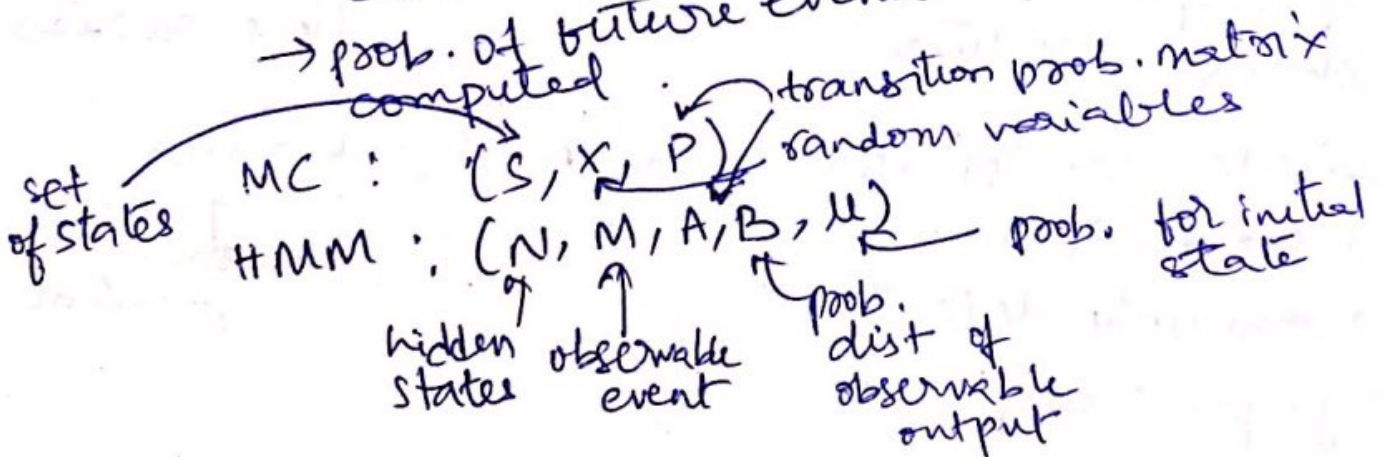→ ORC (outlier removal clustering → uses kmeans)

# Feature selection

using correlation                              using P-value
                                               (not
                                               recommended),
? does it only
  depict linear dependency?
  or a more complicated
  function also?

## Markov chain models
→ history of prev event is known.
→ prob. of transition from 1 event to
  another can be measured
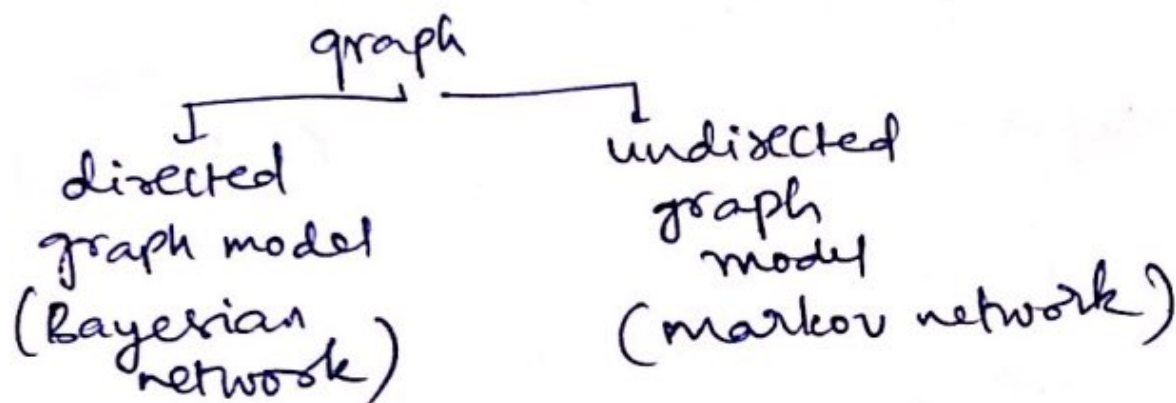→ prob. of future events can be
  computed.

                                    transition prob. matrix
set                         random variables
of states    MC :   (S, X, P)
             HMM :   (N, M, A, B, μ)  — prob. for initial
                       ↑   ↑   ↑              state
                    hidden observable  prob.
                    states  event     dist of
                                      observable
                                      output

# Natural Language Understanding

• intent recognition
• entity recognition → entity
                           ↗    ↘
                       named    numeric
                       entities entities

## NLP, NLU, NLG

→ bayesian networks
→ maximum entropy
→ conditional random field
→ matrix factorization

graphical
models

graph

directed
graph model
(Bayesian
network)

undirected
graph
model
(markov network)

dynamic graphs

structure
of graph
changes with time

node
attributes
can be
time series

→ node classification → node attribut inference
→ link prediction → recommender systems travelling salesman problem
→ community detection
→ graph classification

Node classification in homogeneous graph
→ page rank, centrality measures : baseline models

manual feature engineering to augment
vocabulary-based feature vectors.
with graph-related node features.

* Methods to calculate quantitative values of structural position of a node in a graph.

centrality measures
examples
→ degree centrality
→ betweeness centrality
→ closeness centrality
→ eigenvector centralcty

GNN → utilize static relationships in traing NN on graphs (eg: GCN)

1. GCNN layer → $Z = \sigma ( A' F w + b )$ ← node features matrix

output | activti layer | graph structure | trainable parameters

(graph adjacency matrix)

Markov chain monte Carlo:
→ class of algos for systematic random sampling from high-dimensional probability distribns
→ drawing samples where next sample is dependent on prev sample.

Methods to estimate prediction interval

→ ensemble ANN
→ bayesian method ⎤ will be cooked into
→ monte carlo method ⎦ deeper later.
→ bootstrap method
→ LUBE method (predicty lower & upper bounds)