

# Theory

## Learning From Data

---

---

---

---



# Map

## Paradigms:

supervised  
unsupervised  
reinforcement  
active  
online

## Theory

VC  
bias - variance  
complexity (computation, P, NP, ...)  
bayesian

## Techniques

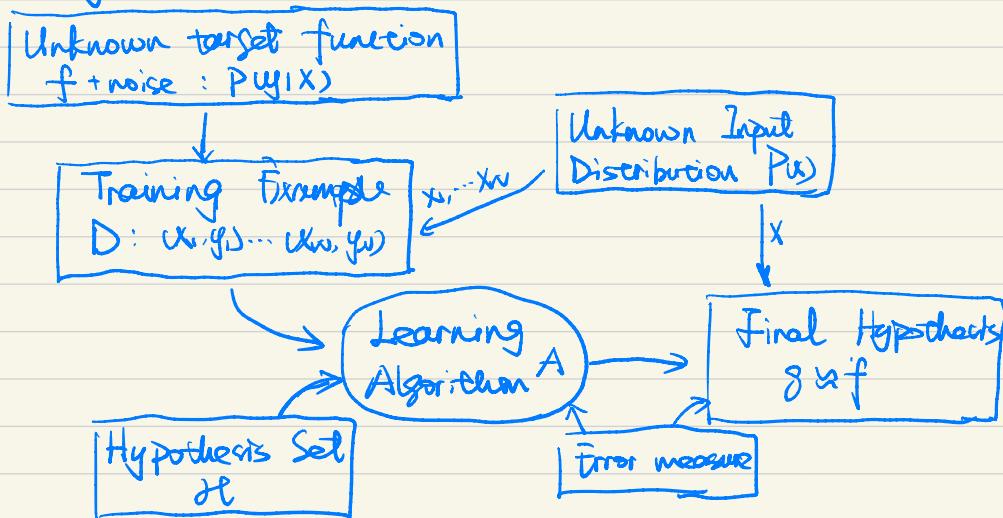
### Model

linear  
neutral networks  
SVM  
nearest neighbors  
RBF  
Gaussian Process  
SVD  
graphical model

### Method

regularization  
validation  
aggregation  
input processing

# Learning Problem



Learning model = hypothesis set + Learning Algorithm

Views

- Statistic: rigorous proof, idealized model in great detail
- ML: less assumption, weaker result, broader application
- Data Mining: More data analysis than prediction

## Learning Feasibility

- infer outside  $D$  with  $D$  in a probabilistic way
- ① A pattern exists
  - ② Can't pin it down mathematically
  - ③ We have data

## Hoeffding's Inequality

Let  $x_1, \dots, x_n$  be independent R.V bounded  $[a_i, b_i]$  a.s.

$$S_n = x_1 + \dots + x_n$$

$$\Rightarrow \begin{cases} P(S_n - E(S_n) \geq t) \leq \exp\left(-\frac{2t^2}{\sum(b_i - a_i)}\right) \\ P(|S_n - E(S_n)| \geq t) \leq 2\exp\left(-\frac{2t^2}{\sum(b_i - a_i)}\right) \end{cases}$$

In sample error

$$\hat{E}_{in}(h) = \frac{1}{N} \sum_i [h(x_i) \neq f(x_i)]$$

= fraction of  $D$  that  $f$  disagree with  $h$

Out sample error

$$E_{out}(h) = P(h(x) \neq g(x)), P \text{ is based on } x$$

$$\downarrow 0 \leq \hat{E} \leq 1$$

$$P(|\hat{E}_{in}(h) - E_{out}(h)| > \epsilon) \leq 2e^{-2\epsilon^2 N}$$

$$\downarrow \mathcal{H} = \{h_1, \dots, h_m\} \quad \xrightarrow{\text{choose } h \text{ based on } \hat{E}_{in}}$$

$$P(|\hat{E}_{in}(h) - E_{out}(h)| > \epsilon) \leq \sum_{h \in \mathcal{H}} P(|\hat{E}_{in}(h) - E_{out}(h)| > \epsilon)$$

$$(2M e^{-2\epsilon^2 N}) \propto \text{complexity of } h$$

Feasibility Question split into 2

1. Can we make sure  $E_{out}(g)$  is close to  $\hat{E}_{in}(g)$   
Hoeffding's Inequality

2. Can we make  $\hat{E}_{in}(g)$  small enough  
Learning algorithm A

Error Measure

- In sample Performance over entire input space  $x$

$$\hat{E}_{in}(h) = \frac{1}{N} \sum_i e(h(x_i), f(x_i))$$

- Out of sample Performance in  $D$

$$E_{out}(h) = \mathbb{E}_x [e(h(x), f(x))]$$

In a noisy set

$p_{ij}(x)$  is what we are looking for

$p_{ij}$  is the relative importance of  $x$  in gauging performance

## VC Dimension

Testing:

$$P(|\hat{E}_{in} - E_{out}| > \epsilon) \leq 2e^{-2\epsilon^2 N}$$

Training:

$$P(|\hat{E}_{in} - E_{out}| > \epsilon) \leq 2M e^{-2\epsilon^2 N}$$

## Effective # of Hypothesis

1. Growth Function group  $h$  based on its behavior on  $D$

- A hypothesis  $h: X \rightarrow \{-1, 1\}$  all input space
- A dichotomy  $h: \{x_1, \dots, x_N\} \rightarrow \{-1, 1\}$ 
  - ↓ Number of Hypotheses  $|H|$  can be infinite
  - Number of dichotomy  $|D(x_1, \dots, x_N)|$  is finite

$$|D(x_1, \dots, x_N)| \leq 2^N$$

Growth function counts the most dichotomy of  $x \in X$

$$m_H(x_1, \dots, x_N) = \max_{x_1, \dots, x_N} |D(x_1, \dots, x_N)| \leq 2^N$$

2. Bound of Growth Function (with break point)

Break Point: If no data set of size  $k$  can be shattered by  $\mathcal{H}$ , then  $k$  is a break point

Bound: If  $m_H(N) < 2^N$  for some  $N$

then  $m_H(N) \leq \frac{K}{\delta} \binom{N}{k}$   $\rightarrow$  polynomial in  $N$  of degree  $k-1$

VC dim (Vapnik-Chervonenkis) of a hypothesis set  $\mathcal{H}$

$vc(\mathcal{H})$  is the largest value of  $N$  for which  $m_H(N) = 2^N$

$$m_H(N) \leq \frac{K}{\delta} \binom{N}{k}$$

The most points that  $\mathcal{H}$  can shatter

## VC dimension

$\text{dvc} \geq N \Leftrightarrow$  there exists  $H$  of size  $N$  such that  $H$  shatters  $D$

$\text{dvc} < N \Leftrightarrow$  No set of  $N$  points can be shattered by  $H$

$\curvearrowleft$  degree of freedom: effective number of hypothesis

3 Replace  $M$  with  $m_{\text{vc}}(w)$

$$P[|\hat{E}_{\text{in}}(g) - E_{\text{out}}(g)| > \varepsilon] \leq 4m_{\text{vc}}(w)e^{-\frac{\varepsilon^2 N}{8}}$$

$\star$  VC generalization bound:  $E_{\text{out}}(g) \leq \hat{E}_{\text{in}}(g) + \sqrt{\frac{8}{N} \ln(\frac{4m_{\text{vc}}(w)}{\varepsilon})}$  with  $1-\delta$

proof: estimate (Replace)  $|\hat{E}_{\text{in}}(g) - E_{\text{in}}(g)| > \varepsilon$

with an artificial event  $|\hat{E}_{\text{in}}(g) - E_{\text{in}}(g)| > \varepsilon$

- A loose bound, but tends to be equally loose hence useful for comparing generalization performance of models

$$\begin{aligned} E_{\text{out}}(g) &\leq \hat{E}_{\text{in}}(g) + \Omega(N, H, S) \\ \Omega(N, H, S) &= \sqrt{\frac{8}{N} \ln\left(\frac{4m_{\text{vc}}(w)}{\varepsilon}\right)} \\ &\leq \sqrt{\frac{8}{N} \ln\left(\frac{4(m_{\text{vc}})^{\text{dvc}} + 1}{\varepsilon}\right)} \end{aligned}$$

- A loose estimate of  $E_{\text{out}}(g)$  as guideline  
Test sample  $\rightarrow$  estimate  $E_{\text{out}}(g)$
- Can extend to other Error Measure

## Bias vs. Variance

### Approximation - Generalization tradeoff

Goal: small  $E_{\text{out}}$ , good approximation of  $f$  out of sample  
complex  $g$  } better chance of approximating  $f$   
| worse chance of generalizing out of sample  
(Have to navigate in  $\mathcal{H}$  through  $D$  for  $g$ )

Bias vs Variance: decomposing  $E_{\text{out}}$  into:

1. How well  $\hat{g}$  can approximate  $f$   
(overall, not in current sample)
2. How well can zoom in on a good  $g \in \mathcal{H}$

Apply to real value  $f$  and square error

$$E_{\text{out}}(g^{(D)}) = E_{\mathcal{D}}[(g^{(D)}(\omega) - f(\omega))^2]$$

↓ remove dependence on  $D$  by expectation

$$\begin{aligned} E_D[E_{\text{out}}(g^{(D)})] &= E_D[E_{\mathcal{D}}[(g^{(D)}(\omega) - f(\omega))^2]] \\ &= E_{\mathcal{D}}[E_D[(g^{(D)}(\omega) - f(\omega))^2]] \end{aligned}$$

↓ average hypothesis

$\bar{g}(\omega) = E_D[g^{(D)}(\omega)] \approx \frac{1}{k} \sum_i^k g^{(D_i)}(\omega)$  on many data set  $D_1 \dots D_k$   
the best can get from the  $\mathcal{H}$  (infinite data)

$$E_D[(g^{(D)}(\omega) - f(\omega))^2]$$

$$= \underbrace{E[(g^{(D)}(\omega) - \bar{g}(\omega))^2]}_{\text{variance}} + (\bar{g}(\omega) - f(\omega))^2$$

bias

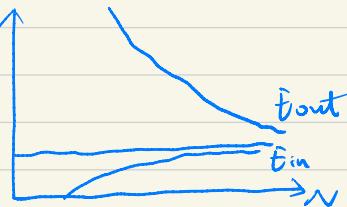
$$\begin{aligned} E_D[E_{\text{out}}(g^{(D)})] &= E_{\mathcal{D}}[\text{bias}(\omega) + \text{var}(\omega)] \\ &= \text{bias} + \text{var} \end{aligned}$$

# Learning Curve

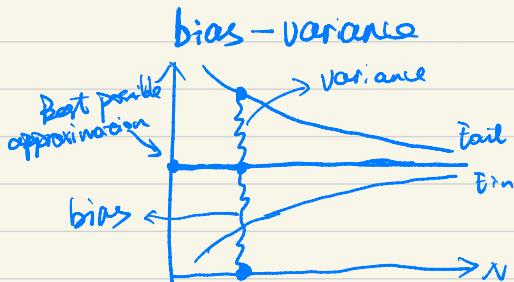
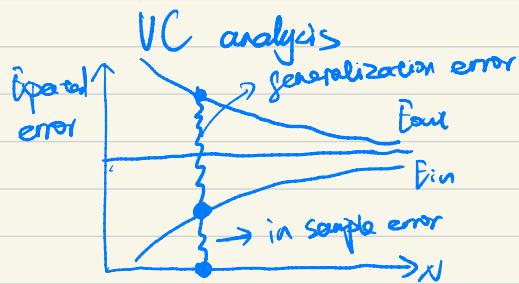
## Simple Model



## Complex Model



Match model complexity with data source  
not target complexity



## Linear Case

$$y = \mathbf{w}^* \mathbf{x} + \text{noise}$$

- ① Best approximation error  $\sigma^2 = \text{Var}(\text{noise})$
- ② Expected in sample error  $\sigma^2(1 - \frac{d+1}{N})$
- ③ Expected out of sample error  $\sigma^2(1 + \frac{d+1}{N})$
- ④ Expected generalization error  $2\sigma^2 \frac{d+1}{N}$

# Overfitting

Overfitting: Ein ↓ Eout ↑

fitting the data more than wanted

(fitting the noise)

$$y = f(x) + \underbrace{\epsilon(x)}_{\sigma^2} \stackrel{\text{def}}{=} \sum_i \alpha_i x_i + \epsilon(x)$$

Overfitting: noise level:  $\sigma^2$  + stochastic noise  
 target complexity:  $D_f$  + deterministic noise  
 data set size:  $N$  -

Deterministic Noise: The part of  $f$  that we can't capture

Noise: something can't capture / deal with

Difference with stochastic noise.

1. depend on  $x$  (complex  $x \rightarrow$  smaller)
2. fixed for a given  $x$

## Noise and bias-variance

$$\mathbb{E}_\theta[(g(w) - f(w))^2] = \mathbb{E}_\theta[(g(w) - \bar{g}(w))^2] + [\bar{g}(w) - f(w)]^2$$

↓ if  $f(w)$  is noisy  $y = f(w) + \epsilon(w)$ ,  $\mathbb{E}[\epsilon] = 0$

$$\mathbb{E}_\theta[(g(w) - y)^2]$$

$$= \mathbb{E}_\theta[(g(w) - \bar{g}(w))^2 + (\bar{g}(w) - f(w))^2 + (\epsilon(w))^2 + \text{cross term}]$$

↓ as  $\bar{f}(w) = 0$

$$\underbrace{\mathbb{E}_\theta[(g(w) - \bar{g}(w))^2]}_{\text{var}} + \underbrace{\mathbb{E}_\theta[(\bar{g}(w) - f(w))^2]}_{\text{bias}} + \underbrace{\mathbb{E}_\theta[(\epsilon(w))^2]}_{\sigma^2}$$

impacted by ↙  
 both noise, capturing a  
 model's susceptibility to  
 being led astray by the noise

where overfit come in ← try to fit noise  
 drag  $g$  away

deterministic noise      stochastic noise  
 Not change with  $N$

## Regularization

Constrain the learning algorithm to improve out-of-sample error  
Necessary if model is too complicated for data amount/quality

### View from VC

$$E_{out}(h) \leq E_{in}(h) + \Omega(h) \quad \text{for all } h \in \mathcal{H}$$

Connect a measure  $\Omega(h)$  for the complexity of an individual hypothesis  
should fit better using simple  $h$  from  $\mathcal{H}$

$$\min_h E_{in}(h) \Rightarrow \min_h [E_{in}(h) + \Omega(h)]$$

### View from bias-var

sacrifice bias for smaller var

### Unconstrained solution

$$\min_h E_{in}(h) \rightarrow W_{un} = (X^T X)^{-1} X^T y$$

Soft-order constraint / weight decay

$$\textcircled{1} \quad \min_h E_{in}(h)$$

$$\text{s.t. } w^T w \leq C \quad (\uparrow \Leftrightarrow \downarrow)$$

$\uparrow$  KKT condition

$$\textcircled{2} \quad \min_h E_{in}(h) + \frac{\lambda}{2} w^T w \quad \text{(augmented error) penalize large weight}$$

$$\rightarrow W_{reg} = (X^T X + \lambda I)^{-1} X^T y$$

### Augmented Error

$$\tilde{E}_{in}(h) = E_{in}(h) + \lambda \underbrace{\Omega(h)}_{\uparrow} \quad \text{regularizer}$$

$$\tilde{E}_{out}(h) \leq E_{in}(h) + \Omega(h)$$

$E_{out}$  is a better proxy of  $E_{out}$  than  $E_{in}$

Regularization is necessary since learning is sensitive to both stochastic and deterministic noise

Best way to constrain learning is in the direction of target

} stochastic noise : high frequency / non-smooth  
} deterministic noise : non-smooth



Usually constrain learning toward smooth h hurt more  
→ overfit noise than to hurt fit useful info

Regularizers:

1. weight decay
2. weight elimination

$$L(w) = \sum_{ijl} \frac{(w_{ijl}^{(0)})^2}{P + (w_{ijl}^{(0)})^2}$$

3. Early stop

choose hyper  $\lambda$ , stop time  
by validation

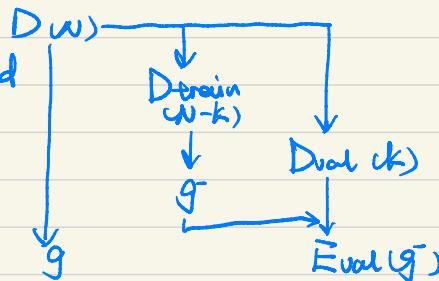
## Validation

Both regularization and validation attempt to minimize  $\hat{E}_{out}$

$$\hat{E}_{out} = \hat{E}_{in} + \underbrace{\text{overfit penalty}}_{\substack{\text{validation estimate} \\ \text{thus}}} \quad \underbrace{\text{regularization estimate this}}$$

## Estimate $\hat{E}_{out}$ with Validation

No decision involved  
just estimation

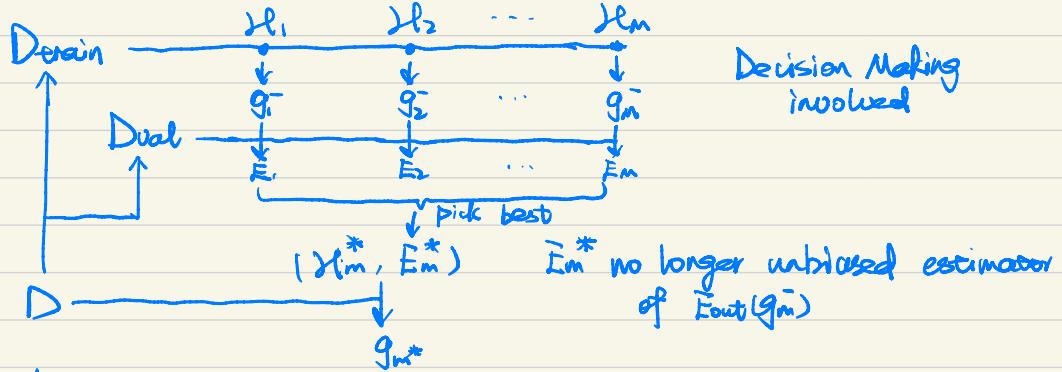


$$\begin{aligned}\hat{E}_{out}(g) &\leq \hat{E}_{in}(g^-) \\ &\leq \hat{E}_{in}(g^-) + O(\frac{1}{k})\end{aligned}$$

$$\hat{E}_{in}(g^-) = \frac{1}{k} \sum_{n \in \text{Dual}} e(g^-(x_n), y_n)$$

- $\hat{E}_{in}(g^-) = \hat{E}_{out}(g^-)$        $e$  depend only on  $x_n$
- $\hat{E}_{in}(e) = \hat{E}_u(e) = \hat{E}_{out}(g^-)$
- $\text{Var}(\hat{E}_{in}(u)) = \frac{s^2}{k}$
- $\hat{E}_{in}(u) = \hat{E}_{out}(u) \pm O(\sqrt{\frac{1}{k}})$
- $k \ll n$ : hard to track  $\hat{E}_{out}(g^-)$  with  $\hat{E}_{in}(g^-)$
- $k \gg n$ :  $\hat{g}^-$  is worse (less data point for training)

## Model selection with validation



↓ Hoeffding

$$\tilde{E}_{out}(g_{m^*}) \leq \tilde{E}_{val}(g_{m^*}) + O(\sqrt{\frac{t_m}{K}})$$

bound is tighter than all Dtrain since  $\lambda_{val} = \{g_1, \dots, g_m\}$  is smaller

- | Dtrain : totally contaminated
- | Dual : slightly contaminated
- | Dtest : totally clean

$$\tilde{E}_{out}(g_{m^*}) \leq \tilde{E}_{out}(g_{m^*}) \leq \tilde{E}_{val}(g_{m^*}) + O(\sqrt{\frac{t_m}{K}})$$

## Cross Validation

$$\tilde{E}_{out}(g) \leq \tilde{E}_{out}(g) \leq \tilde{E}_{val}(g)$$

small k                  large k

leave-one-out cross validation

- ① take data point  $(x_n, y_n)$  away and train with rest  $\rightarrow g_n$
- ②  $e_n = \tilde{E}_{val}(g_n) = e(g_n(x_n), y_n)$
- ③  $E_w = \frac{1}{N} \sum e_n$

Why a good estimator?

$$\text{Define } \tilde{E}_{out}(N) \triangleq \tilde{E}_D[\tilde{E}_{out}(g)]$$

$$\hat{E}_D[e_n] = \hat{E}_{D_n} \left[ \hat{E}_{\text{out}(y_n)} [e(g_n(x_n), y_n)] \right]$$

$$= \hat{E}_{D_n} [\hat{E}_{\text{out}(g_n)}]$$

$$= \hat{E}_{\text{out}(W^{-1})}$$



Cross validation can be used for model selection  
 For computational burden leave one out  $\rightarrow k$ -fold

## Theory vs Practice

Learning theory is guideline but not conclusive proof.

1. noise (stochastic + deterministic) leads to overfitting
2. Regularization helps to prevent overfitting by constraining the model, reducing the impact of the noise while still giving flexibility to fit the data.
3. Validation/CV are useful for estimating  $E_{\text{out}}$ . One important use is model selection

## Kernel Method

Given 2 points  $x$  and  $x' \in \mathcal{X}$

$$k(x, x') = z^T z' \text{ for } z = \phi(x) \text{ in higher dimension}$$

Compute  $k(x, x')$  without transforming  $x$  and  $x'$

### Polynomial kernel

$x \in \mathbb{R}^d$ ,  $\phi: x \rightarrow z$  polynomial of order  $Q$

$$k(x, x') = (1 + x^T x')^Q$$

stay here  
↑

$$= (1 + x_1 x'_1 + x_2 x'_2 + \dots + x_d x'_d)^Q$$

$$= z^T z'$$

too many terms

As long as  $k(x, x')$  is inner product in some  $Z$  space  $\rightarrow$  good

### Mercer's condition

$k(x, x')$  is a valid kernel iff

1. It is symmetric

2. Matrix  $[k(x_i, x_j)]_{ij}$  is positive semi-definite  
for any  $x_1, \dots, x_n$

## Radial Basis Function

Each  $(x_n, y_n) \in X$  influences  $h(x)$  based on  $\|x - x_n\|$  radial

$$h(x) = \sum_{n=1}^N w_n \underbrace{\exp(-r \|x - x_n\|^2)}_{\substack{\text{Basis function} \\ \text{contribution of } n \text{ on } h(x)}} + b$$

Solve for  $w_n$

When  $\hat{x}_n = 0$ :  $h(\hat{x}_n) = y_n$  for  $n=1, \dots, N$

$$\sum_{m=1}^N w_m \exp(-r \|x_m - \hat{x}_n\|^2) = y_n \quad \text{for } n=1, \dots, N$$

$$\begin{bmatrix} \exp(-r \|x_1 - \hat{x}_n\|^2) & \dots & \exp(-r \|x_1 - \hat{x}_N\|^2) \\ \vdots & \vdots & \vdots \\ \exp(-r \|x_N - \hat{x}_n\|^2) & \dots & \exp(-r \|x_N - \hat{x}_N\|^2) \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_N \end{bmatrix} = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}$$

RBF with k centers

$N \gg k$ :  $u_1, \dots, u_k$  instead of  $x_1, \dots, x_N$

$$h(x) = \sum_{k=1}^K w_k \exp(-r \|x - u_k\|^2)$$

$\Downarrow$  estimate  $w_k$  with k-mean, not contaminated

$$\begin{bmatrix} \exp(-r \|x_1 - u_1\|^2) & \dots & \exp(-r \|x_1 - u_K\|^2) \\ \vdots & \vdots & \vdots \\ \exp(-r \|x_N - u_1\|^2) & \dots & \exp(-r \|x_N - u_K\|^2) \end{bmatrix} \begin{bmatrix} w_1 \\ \vdots \\ w_K \end{bmatrix} \approx \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}$$

$$w = (\Phi^T \Phi)^{-1} \Phi^T y$$

RBF can be derived purely on regularization

$$\sum_{n=1}^N (h(\hat{x}_n) - y_n)^2 + \lambda \int_{-\infty}^{\infty} \left( \frac{d^k h}{dx^k} \right)^2 dx$$

Smoothest interpolation

## Learning Principles

### Ovam's Razor

The simplest model that fits the data is also the most plausible \*

#### 1. What's a simple model?

- Complexity of  $h$ : Minimum description length (MDL), order of polynomial
- Complexity of  $H$ : Entropy, VC dimension

#### Link between 2. complexity

$l$  bits specify  $f \rightarrow h$  is one of  $2^l$  of a set  $H$

exception: Looks complex but is one of few - SVM

#### 2. Why is simpler better

better: better out-of-sample performance

- There're fewer simple hypotheses than complex ones  
     $\Downarrow$   $m_{\text{exp}}(H)$

- Less likely to fit a given data set  
     $\Downarrow$   $m_{\text{exp}}(H) / 2^n$

- More significance when it happens

### Sampling Bias

If data is sampled in a biased way. Learning will produce a biased outcome \*

### Data Snooping

If a data set has affected any step in the learning process, its ability to assert the outcome has been compromised \*

most common trap for practitioners

When reusing data set: Trying one model after another on the same data set, you will eventually succeed

- 1. avoid data snooping: strict discipline

- 2. account for data snooping: how much data contamination



# Application

---

---

---

---



# High Dimensional Data

Embedding: a relatively low-dimensional space into  
 $(\mathbb{R}^n)$  which you can translate high dimensional  
vectors

---

---

---

---





# Network Analysis

---

---

---

---



# Spectral Clustering

- Spanning Tree: connect all nodes while minimizing sum of edges
- Steiner Tree: connect selected nodes while minimizing edges
  - NP-problem: need to decide which nodes to include
  - Polynomial Approximation: with spanning tree over those selected nodes
- Community Detection: detect subset of nodes that're more densely connected
- Clustering: determine subsets of points that are 'close' to each other given a similarity measure

Quality function: function that assigns a number to partition

$$\text{Modularity: } Q = \frac{1}{2m} \sum_{i,j} (A_{ij} - P_{ij}) S(i, j)$$

- ①  $P_{ij}$ : expected edge density between  $i$  and  $j$  under null
- ② Erdos-Renyi:  $P_{ij} = \sum_m / m(m-1)$
- ③ Configuration:  $P_{ij} = k_i k_j / (2m-1)$

Modularity maximization: NP

- Louvain Method

Ideal Cluster: identifying nodes that share edges more among themselves compared to outside of themselves.

$$S_{i,j} = \begin{cases} 1 & : \text{group 1} \\ -1 & : \text{group 2} \end{cases} \Rightarrow C = \min_{S \in \{-1, 1\}^n} \sum_{i,j} A_{ij} (1 - S_i S_j)$$

$$n_1, n_2 \text{ size of 2 cluster} \quad \text{s.t.} \quad \sum S_k = n_1 - n_2$$

## Relaxation of clustering

- ① remove integer constraint for s

$$\hat{C}_1 = \min_{X \in \mathbb{R}^{n \times k}} X^T L X, \|X\|_F = 1$$

$$L = D - A = \text{adjacent matrix} - \text{degree matrix}$$

↓  
spectral clustering

# Time Series

---

---

---

---

---



# Time Series Data

Time Series: a collection of realizations of distinct variable at equally spaced points in time.

$$R.V.: X_1 \cdots X_t \cdots X_n \quad (\Omega, P)$$

$$x_i = X_1(w) \cdots x_t = X_t(w) \cdots x_n = X_n(w)$$

- ① time series observations 属于 stochastic process To single realiz
- ② dependence 属于 predict, 但也会有 To statistical info
- ③ time interval 不是 deterministic To

## Deterministic Dependence

常见形式 f mean  $m_x(t) = E(X_t)$

- ① trend:  $m_x(t)$
- ② seasonality:  $s_x(t)$

## Basic Stat

- Marginal Mean:  $m_x(t) = E(X_t)$
- Marginal Variance:  $\text{Var}(X_t) = E[(X_t - m_x(t))^2]$
- Autocovariance  
 $\gamma_{X(S,T)} = \text{cov}(X_S, X_T) = E[(X_S - m_{X(S)}) (X_T - m_{X(T)})]$   
若无相关性, R.V. 为 独立 realization  
↓ assumption

stationary

weak version:

$$m_x(t) = m_x, \text{Var}(X_t) = \sigma_x^2, \text{cov}(X_S, X_T) = \gamma_{X(S-T)}$$

strong version:

distribution  $X_{t+h}, \dots X_{t+nh}$  same for all  $t, n, h$

## Stochastic Dependence

$P_{X_t | X_S}$ ,

$\gamma_{X(S,T)}$  describe linear stochastic dependence

# Stochastic Dependence

$P_{X_t|X_s}$

## Why we need stationary

1. same distribution / stats so we can estimate from time series observations
2. Generalization of CLT / LLN when RV with large gap are kind of independent
3. Rules out trend, seasonality

$$X_t = T_t + S_t + \underset{\text{stationary}}{Y_t}$$

## How to get stationary

1. Remove linear regression on  $t \rightarrow$  remove linear trend
2. Remove nonparametric regression on  $t$   
i.e.: kernel smoothing, polynomial, k-neighbor  
 $\rightarrow$  remove trend / seasonality
3. Remove periodic regression on  $t \rightarrow$  remove seasonality
4. Differentiation  
1st  $\rightarrow$  remove linear trend / expanding variance  
2nd  $\rightarrow$  remove polynomial trend / expanding variance  
other interval  $\rightarrow$  remove seasonality  
...
5. Transformation  
 $\log \rightarrow$  remove expanding variance

## Empirical Estimator

$$\hat{\mu} = \bar{x}_n = \frac{1}{n} \sum_{t=1}^n x_t$$

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{t=1}^n (x_t - \hat{\mu})^2$$

$$\hat{f}(h) = \frac{1}{n} \sum_{t=1}^{n-h} (x_t - \hat{\mu})(x_{t+h} - \hat{\mu}) \text{ for } 1 \leq h < n$$

$$\star \text{Var}(\hat{\mu}) = \frac{1}{n} \sigma^2 + \underbrace{\frac{2}{n^2} \sum_{h=1}^{n-1} (n-h)r(h)}$$

In order for  $\text{Var}(\hat{\mu}) \rightarrow 0$ , the second term need to decay sufficiently fast.  $\Rightarrow$  estimate marginal mean from single path of realization

} stationary  
dependence decay fast

$\Rightarrow$  can estimate expectation from mean in one realization path

# Models

## Autoregressive Process

$$AR(p): X_t = \phi_1 X_{t-1} + \phi_2 X_{t-2} + \dots + \phi_p X_{t-p} + w_t$$

- ACF decay exponentially and never 0
- stationary under conditions

Random walk

$$\begin{cases} \text{No drift: } X_t = X_{t-1} + w_t = X_0 + w_1 + \dots + w_t \\ \text{With drift: } Y_t = S_t + X_t \end{cases}$$

$$= S + Y_{t-1} + w_t$$

$$= S + Y_0 + w_1 + \dots + w_t$$

$$\text{stat, } \mu_{X(t)} = E(X_t) = E(X_0)$$

$$\sigma_{X(t)}^2 = \text{Var}(X_t) = \text{Var}(X_0) + t\sigma_w^2$$

$$\rho_X(s, t) = \text{cov}(X_s, X_t) = \text{Var}(X_0) + \min(s, t)\sigma_w^2$$

$$\mu_Y(t) = S_t + E(Y_0)$$

$$\sigma_Y^2(t) = \text{Var}(Y_0) + t\sigma_w^2$$

$$\rho_Y(s, t) = \text{Var}(Y_0) + \min(s, t)\sigma_w^2$$

## Moving Average

$$MA(q): X_t = w_t + \theta_1 w_{t-1} + \theta_2 w_{t-2} + \dots + \theta_q w_{t-q} \\ = \sum_{h=0}^{q-1} \theta_h w_{t-h}$$

$$\rho(h) = \sum_{j=0}^{q-h} \theta_j \theta_{j+h} \sigma_w^2 \text{ for } 0 \leq h \leq q \Rightarrow \text{stationary}$$

## ARMA/ARIMA

$$\text{ARMA}(p, q): X_t = \phi_1 X_{t-1} + \phi_2 X_{t-2} + \dots + \phi_p X_{t-p} \\ + w_t + \theta_1 w_{t-1} + \theta_2 w_{t-2} + \theta_q w_{t-q}$$

ARIMA: ARMA with differencing

# Fit TS Regression

## Overview

1. transform to stationary
  - log-transform
  - remove trend / seasonalize
  - differentiate successively
2. check for white noise (ACF)
3. If stationary
  - ① finite lag  $\Rightarrow MA(q)$ 
    - compute ACF to get order
    - estimate with ML
    - check residual
  - ② Otherwise  $\Rightarrow AR(p)$ 
    - compute PACF to get order
    - estimate  $\phi_k$  and residual variance  $\hat{\sigma}^2$  via Yule-Walker
    - check residual
  - ③ ARMA(p,q)
    - attempt to fit AR and compute residual
    - attempt to fit MA to residual
    - fit ARMA with p,q from previous steps
    - check residual

	ACF	PACF
AR(p)	decays	0 after p
MA(q)	0 after q	decays
ARMA(p,q)	decays	decays

## AR(p) Parameter Estimation

$$AR(p): X_t = \phi_1 X_{t-1} + \phi_2 X_{t-2} + \dots + \phi_p X_{t-p} + W_t$$

- $p+1$  parameters to estimate:  $\phi_1, \dots, \phi_p, \sigma_w^2$

- If  $X_t$  stationary  $\Rightarrow$  estimate with method of moment

$$(Y_{X(0)}, Y_{X(1)}, \dots, Y_{X(p)}) = \tilde{\gamma}(\phi_1, \phi_2, \dots, \phi_p, \sigma_w^2)$$



Yule-Walker equation

$p+1$  unknown and  
 $p+1$  equations

$$(\hat{\phi}_1, \hat{\phi}_2, \dots, \hat{\phi}_p, \hat{\sigma}_w^2) = \tilde{\gamma}^{-1}(\hat{Y}_{X(0)}, \dots, \hat{Y}_{X(p)})$$

$\underbrace{\qquad\qquad\qquad}_{\text{continuous mapping}}$   
 $\Downarrow$   
 Delta theorem  
 $\Downarrow$   
 consistent estimator

## Yule-Walker:

$$\left. \begin{aligned} E(X_t X_{t+1}) &= E[\phi_1 X_t X_{t-1} + \phi_2 X_t X_{t-2} + \dots + \phi_p X_t X_{t-p} + X_t W_t] \\ E(X_{t+1} X_{t+2}) &= E[\phi_1 X_{t+1} X_{t} + \phi_2 X_{t+1} X_{t-2} + \dots + \phi_p X_{t+1} X_{t-p} + X_{t+1} W_t] \\ &\vdots \\ E(X_{t+p} X_t) &= E[\phi_1 X_{t+p} X_{t-1} + \phi_2 X_{t+p} X_{t-2} + \dots + \phi_p X_{t+p} X_{t-p} + X_{t+p} W_t] \end{aligned} \right)$$



$$\left. \begin{aligned} Y_{X(0)} &= \phi_1 Y_{X(1)} + \phi_2 Y_{X(2)} + \dots + \phi_p Y_{X(p)} + \sigma_w^2 \\ Y_{X(1)} &= \phi_1 Y_{X(0)} + \phi_2 Y_{X(1)} + \dots + \phi_p Y_{X(p-1)} \\ &\vdots \\ Y_{X(p)} &= \phi_1 Y_{X(p-1)} + \phi_2 Y_{X(p-2)} + \dots + \phi_p Y_{X(0)} \end{aligned} \right)$$

# Time Series Analysis

## Partial Correlation

✓ Partial correlation of  $X, Y$  given  $Z$

- regress  $X, Y$  on  $Z$

- $P_{XY|Z} = \text{corr}(X - \hat{X}, Y - \hat{Y})$

✓ Partial correlation of  $X_t, X_{t-h}$  for stationary series

- ACF:  $\rho_{x(h)} = \text{corr}(X_h, X_0) = \frac{\mathbb{E}[(X_h - \mathbb{E}X_h)(X_0 - \mathbb{E}X_0)]}{\text{Var}(X_0)}$

- PACF:  $\rho_{x(h)} = \text{corr}(X_h - \hat{X}_h^{\text{linh-1}}, X_0 - \hat{X}_0^{\text{linh-1}})$   
linear projection on  $X_1, \dots, X_{h-1}$

- Frisch-Waugh-Lovell Theorem

$X_t = \phi_1 X_{t-1} + \dots + \phi_h X_{t-h} + \hat{X}_t$  is best linear prediction  
 $\Downarrow$

$$\rho_{x(h)} = \phi_h$$

## Information Criterion

$$\left. \begin{array}{l} AIC(k) = \underbrace{-2 \log\text{-likelihood}}_{\text{model fit}} + 2k \\ BIC(k) = -2 \log\text{-likelihood} + k \ln(n) \end{array} \right\} \begin{array}{l} \text{complexity penalty} \\ \text{BIC} > AIC \end{array}$$

# Gaussian Process

---

---

---

---



# Gaussian Process

$N$ -dim multivariate Gaussian

$$\text{joint} \quad \begin{cases} X = \begin{bmatrix} x_1 & x_1 \in \mathbb{R}^d \\ x_2 & x_2 \in \mathbb{R}^{N-d} \end{bmatrix} \\ U = \begin{bmatrix} u_1 & u_1 \in \mathbb{R}^d \\ u_2 & u_2 \in \mathbb{R}^{N-d} \end{bmatrix} \\ \Sigma = \begin{bmatrix} \Sigma_{11} \in \mathbb{R}^{d \times d} & \Sigma_{12} \in \mathbb{R}^{d \times (N-d)} \\ \Sigma_{21} \in \mathbb{R}^{(N-d) \times d} & \Sigma_{22} \in \mathbb{R}^{(N-d) \times (N-d)} \end{bmatrix} \end{cases}$$

$$\text{conditionally } u_{x_1|x_2} = u_1 + \Sigma_{12} \Sigma_{22}^{-1} (x_2 - u_2)$$

$$\Sigma_{x_1|x_2} = \Sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}$$

$$x_{1|2} \sim N(u_{x_1|x_2}, \Sigma_{x_1|x_2})$$

## Special covariance

observed value:  $y_1, \dots, y_N$

location:  $x_1, \dots, x_N$

$$\text{RBF: } \text{cov}(y_i, y_j) = k(x_i, x_j) = \exp\left(-\frac{\|x_i - x_j\|^2}{2\sigma^2}\right)$$

$$\Sigma_N = \begin{bmatrix} \text{cov}(y_1, y_1) & \cdots & \text{cov}(y_1, y_N) \\ \vdots & \ddots & \vdots \\ \text{cov}(y_N, y_1) & \cdots & \text{cov}(y_N, y_N) \end{bmatrix} = \begin{bmatrix} k(x_1, x_1) & \cdots & k(x_1, x_N) \\ \vdots & \ddots & \vdots \\ k(x_N, x_1) & \cdots & k(x_N, x_N) \end{bmatrix}$$

$$k_*^T = [k(x_*, x_1), \dots, k(x_*, x_N)]$$

$$\text{conditionally } u_{*|1:N} = u_* + k_*^T k_N^{-1} (y_{1:N} - u_{1:N})$$

$$\Sigma_{*|1:N} = \Omega_*^2 - k_*^T k_N^{-1} k_*$$

predict:

$$\hat{y}_* = \hat{f}(x_*) = u_* + \underset{1:N \times N \times 1}{k_*^T k_N^{-1} y_{1:N}}$$

=

$$\begin{cases} \textcircled{1} \sum_{i=1}^N w_i^* y_i & \text{weighted average of } y_{1:N} \\ \textcircled{2} \sum_i k_{*,i} x_i & \text{linear combination of } N \text{ non-linear features, same } d \text{ for all prediction} \end{cases}$$

## Gaussian Process

a collection of random variables, any finite number of which is Gaussian

$$\left\{ \begin{array}{l} m(x) = \mu_x = E[f(x)] \\ k(x, x^T) = \text{cov}(f(x), f(x^T)) \end{array} \right.$$

## Measurement Noise

$$y_i = f(x_i) + \epsilon_i, \quad \epsilon_i \sim N(0, \sigma^2)$$

$$\Rightarrow \left\{ \begin{array}{l} \hat{m}_{*|1:N} = \mu_* + k_*^T (K_N + \sigma^2 I)^{-1} y_{1:N} \\ \hat{\sigma}_{*|1:N}^2 = k_*(x_*, x_*) - k_*^T (K_N + \sigma^2 I)^{-1} k_* \end{array} \right.$$

adjust less due to noise

## Entropy Quantify uncertainty

$$H(Y) = \sum_y P(y) \log \frac{1}{P(y)} = - \sum_y P(y) \log P(y)$$

measure of surprise  
 expected surprise

• Joint Entropy:

$$H(Y, Z) = - \sum_{y,z} P(y, z) \log P(y, z) \leq H(Y) + H(Z)$$

• Conditional Entropy:

$$H(Y|Z) = H(Y, Z) - H(Z)$$

• Mutual Information

$$I(Y, Z) = H(Y) - H(Y|Z) = \sum_z P(z) \log \frac{P(y, z)}{\underbrace{P(y)P(z)}_{\text{if independent}}} \downarrow \text{distance from independence}$$

