

Theory

Learning From Data



Map

Paradigms:

supervised
unsupervised
reinforcement
active
online

Theory

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

Techniques

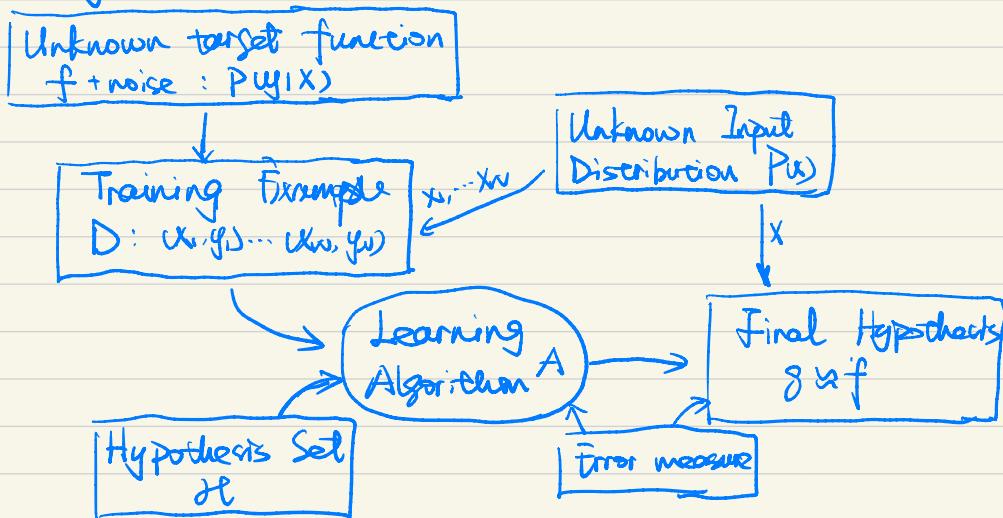
Model

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 $m_H(N) \leq \frac{k!}{(k-N)!} \binom{N}{k} \rightarrow$ 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 $m_H(N) = 2^N$

$$m_H(N) \leq \frac{4e}{\delta} \binom{N}{k}$$

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_{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_{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

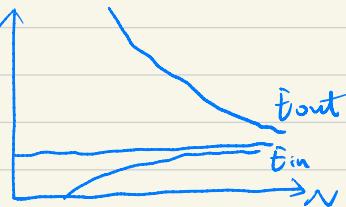
$$\begin{aligned} E_D[E_{\text{out}}(g^{(D)})] &= E_{\mathcal{D}}[\text{bias}(\omega) + \text{var}(\omega)] \\ &= \text{bias} + \text{var} \end{aligned}$$

Learning Curve

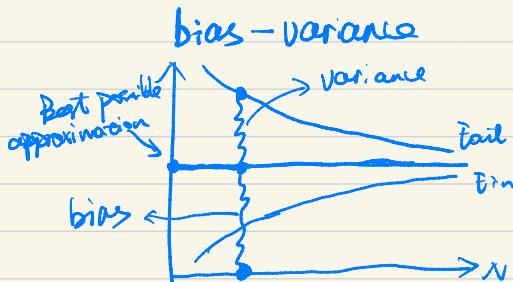
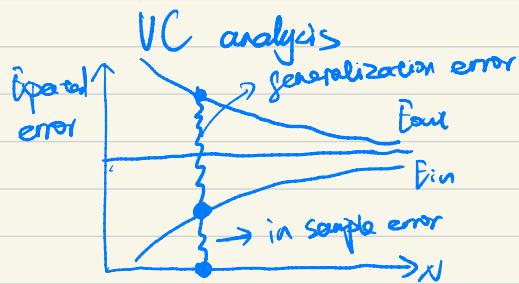
Simple Model



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: σ^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

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

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

E_{aug} 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

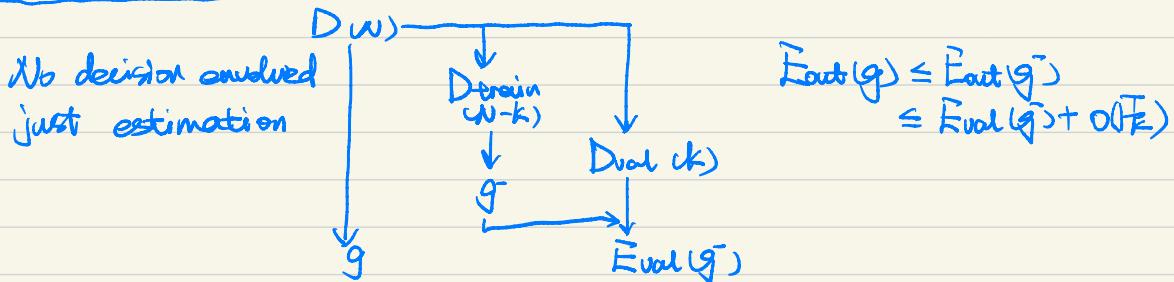
choose hyper λ , stop time
by validation

Validation

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

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

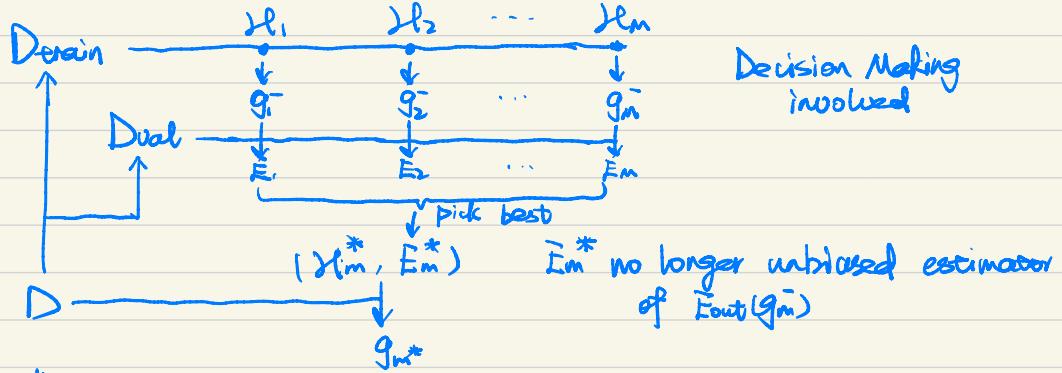
Estimate \hat{E}_{out} with Validation



$$\hat{E}_{\text{val}}(g) = \frac{1}{k} \sum_{x_i \text{ test}} e(g(x_i), y_i)$$

- $\hat{E}_{\text{val}}(\hat{E}_{\text{val}}(g)) = \hat{E}_{\text{in}}(g)$ e depend only on x_i
 $\hat{E}_{\text{val}}(e) = \hat{E}_{\text{in}}(e) = \hat{E}_{\text{out}}(g)$
- $\text{Var}(\hat{E}_{\text{val}}(h)) = \frac{s^2}{k}$
 $\hat{E}_{\text{val}}(h) = \hat{E}_{\text{out}}(h) \pm \sigma \sqrt{\frac{1}{k}}$
- $k \text{ small: hard to track } \hat{E}_{\text{out}}(g)$ with $\hat{E}_{\text{val}}(g)$
- $k \text{ large: } g$ 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 $E_w = \frac{1}{N} \sum e_n$

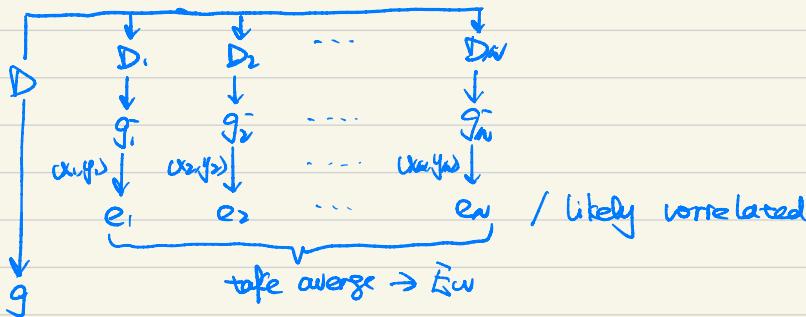
Why a good estimator?

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

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

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

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



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_{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

Discrete Latent Variable

Modeling of distribution over observed/latent variables



K - mean

Objective : $\min_{\mu_1 \dots \mu_K} \sum_{n=1}^N \text{rank} \|x_n - \mu_k\|^2$

Iterative process:

① Initialize μ_k

② Iterate until no change (convergence)

- E: keep μ_k fixed, $\min J$ w.r.t rank
- $$\text{rank} = \begin{cases} 1 & , k = \arg \min_j \|x_n - \mu_j\|^2 \\ 0 & , \text{otherwise} \end{cases}$$

- M: keep rank fixed, $\min J$ w.r.t μ_k

$$\mu_k = \frac{\sum_n \text{rank}_n x_n}{\sum_n \text{rank}_n}$$

EM for Mixture Gaussian

1. Initialize μ_k , Σ_k and π_k

2. E-step

$$\gamma(z_{nk}) = \frac{\pi_k N(x_n | \mu_k, \Sigma_k)}{\sum_j \pi_j N(x_n | \mu_j, \Sigma_j)}$$

3. M-step

$$\mu_k^{\text{new}} = \frac{1}{N_k} \sum \gamma(z_{nk}) x_n$$

$$\Sigma_k^{\text{new}} = \frac{1}{N_k} \sum \gamma(z_{nk}) (x_n - \mu_k^{\text{new}}) (x_n - \mu_k^{\text{new}})^T$$

$$\pi_k^{\text{new}} = \frac{N_k}{N} / N$$

$$N_k = \sum \gamma(z_{nk})$$

4. Check change in log-likelihood and stop if converge.

Mixture Gaussian

① Binary: $P(z_k=1) = \pi_k, \sum \pi_k = 1$

② Conditional Gaussian $P(x|z_k=1) = N(x|\mu_k, \Sigma_k)$

$$\left\{ P(z) = \prod_{k=1}^K \pi_k^{z_k} \right.$$

$$\left. P(x|z) = \prod_{k=1}^K N(x|\mu_k, \Sigma_k)^{z_k} \right.$$

$$P(x) = \sum_z P(z) P(x|z)$$

$$= \sum_{k=1}^K \pi_k N(x|\mu_k, \Sigma_k)$$

$z_k=1$ prior: π_k

$$\text{posterior: } \hat{\pi}(z_k) = P(z_k=1|x) = \frac{\pi_k N(x|\mu_k, \Sigma_k)}{\sum_j \pi_j N(x|\mu_j, \Sigma_j)}$$

Log-likelihood

$$\ln p(X|\pi, \mu, \Sigma) = \sum_{n=1}^N \ln \left(\sum_{k=1}^K \pi_k N(x_n|\mu_k, \Sigma_k) \right)$$

$$X: N \times D \quad Z: N \times K$$

EM in general

General technique to find ML for probabilistic model with latent variables

Goal: Maximize likelihood over all latent variables
 $\text{Max } P(x|\theta) = \underbrace{\sum P(x,z|\theta)}_{\text{discrete } z} = \underbrace{\int P(x,z|\theta) dz}_{\text{continuous } z}$

$\gamma \theta$: parameter

γz : latent variables, x : observed variables

Idea:

optimization over $P(x|\theta)$ is difficult due to latent
but optimization over $P(x,z|\theta)$ is easier

↓ introduce $q(z)$ on latent variable, for any $q(z)$
 $P(x,z|\theta) = P(z|x,\theta) P(x|\theta)$

$$\ln P(x|\theta) = \underbrace{\sum_z q(z) \ln \left(\frac{P(x,z|\theta)}{q(z)} \right)}_{D(q,\theta)} - \underbrace{\sum_z q(z) \ln \left(\frac{P(z|x,\theta)}{q(z)} \right)}_{KL(q||p)}$$

$$= \sum_z q(z) \ln \left(\frac{P(x,z|\theta)}{P(z|x,\theta)} \right)$$

$$= \sum_z q(z) \ln P(x|\theta) = \ln P(x|\theta) \sum_z q(z)$$

$KL(p||q)$: KL Divergence of γ posterior $P(z|x,\theta)$

$$\underbrace{\geq 0}_{\downarrow}$$

$D(q,\theta)$ is a lower bound on $\ln P(x|\theta)$

"=" hold when $KL(p||q) = 0$



$$q(z) = P(z|x,\theta)$$

1. E-step

holding θ^{old} constant, changing $q(z)$ to maximize lower bound $L(q, \theta)$

$\Downarrow P(x|\theta^{\text{old}})$ does not depend on $q(z)$, reach maximum when $\text{KL}(P||q) = 0$

$$q(z) = P(z|x, \theta^{\text{old}})$$

2. M-step

holding $q(z)$ constant, changing θ to maximize lower bound $L(q, \theta)$

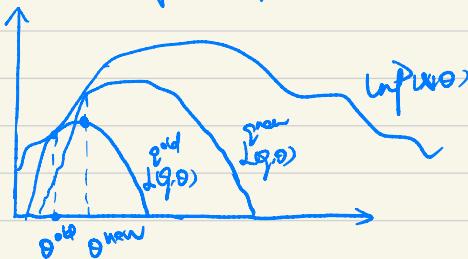
$$\theta^{\text{new}} = \underset{\theta}{\operatorname{argmax}} L(q, \theta)$$

$$\Downarrow q(z) = P(z|x, \theta^{\text{old}})$$

$$L(q, \theta) = \underbrace{\sum z_i P(z_i|x, \theta^{\text{old}}) \ln P(x, z_i|\theta)}_{\text{entropy of } q} - \underbrace{\sum z_i P(z_i|x, \theta^{\text{old}}) \ln P(z_i|x, \theta^{\text{old}})}_{\text{entropy of } P}$$

$$= Q(\theta, \theta^{\text{old}}) + \text{const}$$

expectation of complete data log likelihood



Another Perspective

$$\text{Likelihood } L(\theta) = \log P(x|\theta) = \log \int p(x, z|\theta) dz \\ = \log \int q(z) \frac{p(x, z|\theta)}{q(z)} dz \geq \underbrace{\int q(z) \log \frac{p(x, z|\theta)}{q(z)} dz}_{\text{Jensen's}} \quad \text{Lower Bound}$$

\hat{E} -step

$$q_k(z) = \underset{q(z)}{\operatorname{argmax}} \int q(z) \frac{p(x, z|\theta_{k-1})}{p(z)} dz$$

$$\Rightarrow q_k(z) = p(z|x, \theta_{k-1})$$

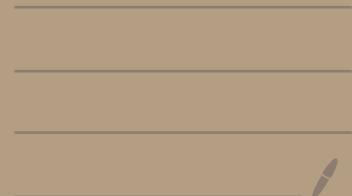
M -step

$$\theta_k = \underset{\theta_k}{\operatorname{argmax}} \int q_k(z) \frac{p(x, z|\theta)}{q(z)} dz$$

\Downarrow

Local Maximum

Continuous Latent
Variables



PCA

Formulation

1. Maximum variance formulation

$$\{x_n\}_{n=1, \dots, N} \quad x_n \in \mathbb{R}^D \quad \bar{x} = \frac{1}{N} \sum x_n$$

$\downarrow u_i^T x_n$ projection as a scalar

$$\text{var: } \frac{1}{N} \sum (u_i^T x_n - u_i^T \bar{x})^2 = u_i^T S u_i, \quad S = \frac{1}{N} \sum (x_n - \bar{x})(x_n - \bar{x})^T$$

For following
optimizing
orthogonal
to previous
transformation
vectors

$$\text{Max } u_i^T S u_i,$$

$$\text{s.t. } u_i^T u_i = 1 \Rightarrow \text{scale}$$

\Downarrow Lagrangian

$$u_i^T S u_i + \lambda_i (1 - u_i^T u_i)$$

\Downarrow first order

$$S u_i = \lambda_i u_i \quad \text{must be an eigenvector}$$

2. Minimum Error formulation

Introducing complete orthonormal basis vector $\{u_i\}_{i=1, \dots, D}$

$$x_n = \sum_{i=1}^D z_{ni} u_i$$

$$\begin{array}{|c|c} \hline \text{coordinate} & (x_{n1}, \dots, x_{nD}) \\ \hline \text{switch} & \downarrow \\ & (z_{n1}, \dots, z_{nD}) \end{array}$$

\Downarrow $x_{nj} = x_n^T u_j$ new coordinate with basis u point

$$\begin{cases} x_n = \sum_{i=1}^D (z_{ni} u_i) u_i \\ x_n = \sum_{i=1}^D z_{ni} u_i + \sum_{i=1}^D b_i u_i \end{cases}$$

$\left\{ \begin{array}{l} z_{ni}: \text{depend on particular data} \\ b_i: \text{same for all} \end{array} \right.$

$$\min_{u, z, b} \frac{1}{N} \sum \|x_n - \hat{x}_n\|^2$$

\Downarrow

$$\frac{\partial}{\partial z} = 0 \Rightarrow z_{nj} = x_n^T u_j$$

$$\frac{\partial}{\partial b} = 0 \Rightarrow b_j = \bar{x}^T u_j$$

$$\hat{x}_n - x_n = \sum_{i=1}^D \{(x_n - \bar{x})^T u_i\} u_i$$

$$\begin{aligned}
 \hat{x}_n &= \sum_{i=1}^m z_{ni} u_i + \sum_{i=1}^D b_i u_i \\
 &\Downarrow \quad x_n = \sum_{i=1}^D (x_n^T u_i) u_i \\
 &\quad \min \|x_n - \hat{x}_n\|^2 \\
 x_n - \hat{x}_n &= \sum_{i=1}^D \{ (x_n - \bar{x})^T u_i \} u_i \\
 &\Downarrow \\
 x_n^T J &= \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^D (x_n^T u_i - \bar{x}^T u_i)^2 = \sum_{i=1}^D u_i^T S u_i \\
 &\Downarrow \\
 S u_i &= \lambda_i u_i, \quad J = \sum_{i=1}^D \lambda_i
 \end{aligned}$$

Application

1. data compression

$$\begin{aligned}
 \hat{x}_n &= \sum_{i=1}^m z_{ni} u_i + \sum_{i=1}^D b_i u_i \\
 &\Downarrow \quad z_{nj} = x_n^T u_j \\
 &\quad b_j = \bar{x}^T u_j \\
 &= \sum_{i=1}^m (x_n^T u_i) u_i + \sum_{i=1}^D (\bar{x}^T u_i) u_i \\
 &= \bar{x} + \sum_{i=1}^m (x_n^T u_i - \bar{x}^T u_i) u_i \quad \text{only consider deviation from mean for first } m
 \end{aligned}$$

2. normalization

$$\begin{aligned}
 S u &= U L \\
 (\dots u_i \dots) &\quad (\dots \lambda_i \dots) \\
 y_n &= L^{-1} U^T (x_n - \bar{x}) \\
 &\quad \left| \begin{array}{l} \bar{y} = 0 \\ \frac{1}{n} \sum y_n y_n^T = I \end{array} \right.
 \end{aligned}$$

3. Dimension Reduction

transform: $\tilde{z}_{n \times k} = X_{n \times p} V_{p \times k}$ first k eigenvectors

decode: $\hat{X}_{n \times p} = \tilde{z}_{n \times k} V_{k \times p}^T$

Probabilistic PCA

$$X = WZ + u + \sigma$$

$$\cdot P(Z) = N(Z|0, I) \quad \text{dim } M$$

$$\cdot P(X|Z) = N(X|WZ + u, \sigma^2 I) \quad \text{dim } D$$

Likelihood: $P(X) = \int p(x|z)p(z)dz = N(X|u, C)$

$$C = \mathbb{E}[XZ] = \mathbb{E}[(WZ + u)(WZ + u)^T]$$

$$= WW^T + \sigma^2 I$$

Posterior: $p(z|x) = N(z|M^{-1}W^T(x-u), \sigma^2 M^{-1})$

$$M = W^T W + \sigma^2 I$$

MLE

$$\ln P(X|u, w, \sigma^2) = \sum_i \ln p(x_i|w, u, \sigma^2)$$

set derivative $\hat{u} = \bar{x}$

$$\hat{w} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T$$

$$\hat{\sigma}^2 = \frac{1}{D-M} \sum_{i=1}^n \lambda_i$$

EM

① Write down complete-data log likelihood

$$\ln P(X, Z|u, w, \sigma^2) = \sum_i \{ \ln p(x_i|z_i, w) + \ln p(z_i) \}$$

② Take expectation with respect to posterior of latent

$$\begin{aligned} \mathbb{E}[\ln p(x_i|z_i, u, w, \sigma^2)] &= -\frac{n}{2} \left\{ \frac{D}{2} \ln(2\pi\sigma^2) + \frac{1}{2} \text{Tr}(E[Z_i Z_i^T]) \right. \\ &\quad + \frac{1}{2\sigma^2} \|x_i - u\|^2 - \frac{1}{2\sigma^2} E(Z_i)^T W^T (x_i - u) \\ &\quad \left. + \frac{1}{2\sigma^2} \text{Tr}(E(Z_i Z_i^T) W^T W) + \frac{M}{2} \ln(2\pi) \right\} \end{aligned}$$

③ Use old poem to evaluate (E-step)

$$\hat{E}(Z_i) = M^{-1} W^T (x_i - \bar{x})$$

$$E(Z_i Z_i^T) = \sigma^2 M^{-1} + E(Z_i) E(Z_i)^T$$

④ Maximize keep posterior fixed (M-step)

$$W_{\text{new}} = \left[\frac{1}{n} \sum_i (x_i - \bar{x}) \hat{E}(Z_i)^T \right] \left[\frac{1}{n} \hat{E}(Z_i) \hat{E}(Z_i)^T \right]$$

$$\sigma_{\text{new}}^2 = \frac{1}{ND} \sum_i \{ \|x_i - \bar{x}\|^2 - 2\hat{E}(Z_i)^T W_{\text{new}} (x_i - \bar{x}) + \text{Tr}(\hat{E}(Z_i Z_i^T) W_{\text{new}}^T W_{\text{new}}) \}$$

Factor Analysis

$$y = \Lambda X + \varepsilon \quad \varepsilon \sim N(0, \Psi) \xrightarrow{\text{diagonal}}$$

factor loading matrix

Non-uniqueness of Λ :

- orthogonal matrix $U : UU^T = I$
- $y = \Lambda X + \varepsilon = \Lambda U U^T X + \varepsilon = (\Lambda U)(U^T X) + \varepsilon$

Connection with PCA:

when $\Psi = \sigma^2 I$, Λ is made of the first k eigenvectors of covariance matrix of y .

EM

E : estimate X with parameters from last step

M : adjust parameters by maximizing the lower bound on the log-likelihood.

Sequential

Data



Sequence Latent State Model

Markov Model

No assumption:

$$P(x_1, \dots, x_N) = P(x_1) \prod_{n=2}^N P(x_n | x_1, \dots, x_{n-1})$$

2-nd order Markov

$$P(x_1, \dots, x_N) = P(x_1) P(x_2 | x_1) \prod_{n=3}^N P(x_n | x_{n-1}, x_{n-2})$$

State Space Model

$$P(x_1, \dots, x_N, z_1, \dots, z_N) = P(z_1) \prod_{n=2}^N P(z_n | z_{n-1}) \prod_{n=1}^N P(x_n | z_n)$$

Sequence latent state models $P(x|z)$

parametric } State-Space model: continuous z
Non-parametric : RNN, LSTM, Bayesian neural networks } Hidden Markov model: discrete z

$$P(x_n | x_{n-1}) \rightarrow P(x_n | z_n)$$

HMM (Hidden markov Model)

$$P(X, Z | \theta) = P(Z_1 | X) \left[\prod_{m=2}^N P(Z_m | Z_{m-1}, A) \right] \prod_{m=1}^M P(X_m | Z_m, \phi)$$

initial state transition conditional probability

$\theta = \{Z, A, \phi\}$

$X = \{X_1, \dots, X_N\}$ observed variable

$Z = \{Z_1, \dots, Z_N\}$ state variable

EM for HMM

① $\hat{\ell}$ -step : find the posterior $P(Z|X, \theta^{old})$

$$k \times 1: r(Z_n) = P(Z_n | X, \theta^{old}) = \alpha(Z_n) \beta(Z_n) / P(X)$$

$$k \times k: \xi(Z_{m-1}, Z_m) = P(Z_{m-1}, Z_m | X) = \alpha(Z_{m-1}) P(X_m | Z_m) P(Z_m | Z_{m-1}) \beta(Z_m) / P(X)$$

$$\alpha(Z_n) = P(X_1, \dots, X_n, Z_n)$$

$$= P(X_n | Z_n) \sum_{z_{n-1}} P(X_1, \dots, X_{n-1}, z_{n-1}) P(Z_n | z_{n-1})$$

$$= P(X_n | Z_n) \sum_{z_{n-1}} \alpha(Z_{n-1}) P(Z_n | z_{n-1})$$

$$\beta(Z_n) = P(X_{n+1}, \dots, X_N | Z_n)$$

$$= \sum_{z_{n+1}} P(X_{n+2}, \dots, X_N | Z_n) P(X_{n+1} | Z_n) P(Z_{n+1} | Z_n)$$

$$= \sum_{z_{n+1}} \beta(Z_{n+1}) P(X_{n+1} | Z_n) P(Z_{n+1} | Z_n)$$

② M -step : Maximize posterior expectation of the log complete likelihood

$$Q(\theta, \theta^{old}) = \sum_i P(Z_i | X, \theta^{old}) \ln P(X, Z | \theta)$$

$$= \sum_k r(Z_k) \ln \pi_k + \sum_h \sum_j \xi(Z_{m-1}, j, Z_m) \ln A_{jk}$$

$$+ \sum_h \xi(Z_k) \ln \beta(X_k | \phi_k)$$

↓ Lagrangian Multiplier

$$\pi_k = \frac{r(Z_k)}{\sum_j r(Z_j)}$$

$$A_{jk} = \frac{\sum_{m=2}^N \xi(Z_{m-1}, j, Z_m)}{\sum_h \xi(Z_{m-1}, h, Z_m)}$$

$$\phi_k \text{ only depend on the last term } \left\{ \begin{array}{l} \mu_k = \frac{\sum r(Z_k) (X_k - \bar{X}) (X_{k-1} - \bar{X})}{\sum r(Z_k)} \\ \Sigma_k = \frac{\sum r(Z_k) (X_k - \bar{X}) (X_{k-1} - \bar{X})}{\sum r(Z_k)} \end{array} \right.$$

State Space Model

$$Z_n = AZ_{n-1} + W_n \quad W \sim N(W|0, P)$$

$$X_n = CZ_n + V_n \quad V \sim N(V|0, \Sigma)$$

$$Z_0 = u_0 + U \quad U \sim N(U|0, P_0)$$

$$\Theta = \{A, P, C, \Sigma, u_0, P_0\}$$

$\hat{E}M$: inference problem: posterior marginals for latent
 $\hat{L}M$: learning problem: update Θ

Inference : \hat{E} step

posterior Marginal	$\hat{P}(Z_n) = P(Z_n X_1, \dots, X_n) = N(Z_n \hat{u}_n, \hat{V}_n)$
-----------------------	---

$$\left. \begin{array}{l} \hat{u}_n = \underbrace{A\hat{u}_{n-1}}_{\text{prediction}} + \underbrace{k_n(X_n - CA\hat{u}_{n-1})}_{\text{correction}} \\ \hat{V}_n = [I - k_n C] P_{n-1} \end{array} \right.$$

$$\begin{aligned} & P(Z_n | X_1, \dots, X_{n-1}) \rightarrow P(Z_n | X_1, \dots, X_{n-1}) \xrightarrow{X_n} P(Z_n | X_1, \dots, X_n) \\ & C_n = P(X_n | X_1, \dots, X_{n-1}) \end{aligned}$$

$$= N(X_n | CA\hat{u}_{n-1}, C\hat{P}_{n-1}C^T + \Sigma)$$

- $P_{n-1} \equiv AV_{n-1}A^T + P$

- kalman Gain: $k_n = P_{n-1}C^T(C\hat{P}_{n-1}C^T + \Sigma)^{-1}$

Marginal	$\gamma(Z_n) = P(Z_n X) / P(X) = \alpha(Z_n) \beta(Z_n) / \alpha(X)$ $= N(Z_n \hat{u}_n, \hat{V}_n)$
----------	---

$$\left. \begin{array}{l} \hat{u}_n = \hat{u}_n + J_n (\hat{u}_{n+1} - A\hat{u}_n) \\ \hat{V}_n = V_n + J_n (\hat{V}_{n+1} - P_n) J_n^T \\ \cdot J_n = V_n A^T (P_n)^{-1} \\ \cdot \text{cov}(Z_{n+1}, Z_n) = J_n \hat{V}_n \end{array} \right.$$

Learning: M step

$$\ln p(x, z | \theta) = \ln p(z_1 | u_0, p_0) + \sum_{n=1}^M \ln p(z_n | z_m, A, T) \\ + \sum_{n=1}^M \ln p(u_n | z_n, C, z)$$

↓ expectation with $P(z|x, \theta^{old})$ posterior distribution

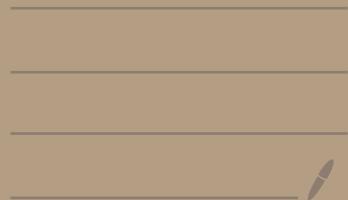
$$Q(\theta, \theta^{old}) = E_{z|x, \theta^{old}} [\ln p(x, z | \theta)]$$

↓ maximize

$$\left\{ \begin{array}{l} u_0^{new} = E(z_1) \\ P_0^{new} = E(z_1 z_1^T) - E(z_1) E(z_1)^T \\ A^{new} = \left(\sum_{n=1}^N E(z_n z_{n-1}^T) \right) \left(\sum_{n=1}^N E(z_n z_n^T) \right)^{-1} \\ \Gamma^{new} = \frac{1}{N-1} \sum_{n=1}^N \{ E(z_n z_n^T) - A^{new} E(z_{n-1} z_{n-1}^T) \\ \quad - E(z_n z_{n-1}^T) (A^{new})^T + A^{new} E(z_n z_n^T) (A^{new})^T \} \\ C^{new} = \left(\sum_{n=1}^N x_n E(z_n^T) \right) \left(\sum_{n=1}^N E(z_n z_n^T) \right)^{-1} \\ \Sigma^{new} = \frac{1}{N} \sum_{n=1}^N x_n x_n^T - (C^{new})^T E(z_n) x_n \\ \quad - x_n E(z_n^T) C^{new} + (C^{new})^T E(z_n z_n^T) C^{new} \end{array} \right.$$

- $E(z_n) = \hat{u}_n$
- $E(z_n z_{n-1}^T) = \hat{v}_n J_{n-1}^T + \hat{u}_n \hat{u}_{n-1}^T$
- $E(z_n z_n^T) = \hat{v}_n + \hat{u}_n \hat{u}_n^T$

High Dimensional Data



- ① PCA: projection that spreads data
- ② Multidimensional scaling: retain original distance
- ③ Stochastic neighbor embedding: non linear method \Rightarrow keep close points together

MDS (multidimensional scaling)

$$Y = \underset{Y}{\operatorname{argmin}} \sum_{i=1}^n \sum_{j=1}^n (D_{ij} - \|y_i - y_j\|^2)$$

$$\left| \begin{array}{l} D_{ij} = \text{dist}(x_i, x_j) \quad x_i \in \mathbb{R}^P \\ y_i \in \mathbb{R}^q \quad q < P \\ \Downarrow D_{ij} = \|x_i - x_j\|_2^2 \end{array} \right.$$

$$\min_Y \text{trace}(XX^T - YY^T)^2$$

$$\Rightarrow Y = U\Lambda, \quad \Lambda \text{ eigenvalue of } XX^T$$

SNE (stochastic neighbor embedding)

- ① probability of picking j as neighbor of i:

$$P_{ij} = \frac{\exp(-D_{ij}^2)}{\sum_{k \neq i} \exp(-D_{ik}^2)}$$

- ② probability in low dimension

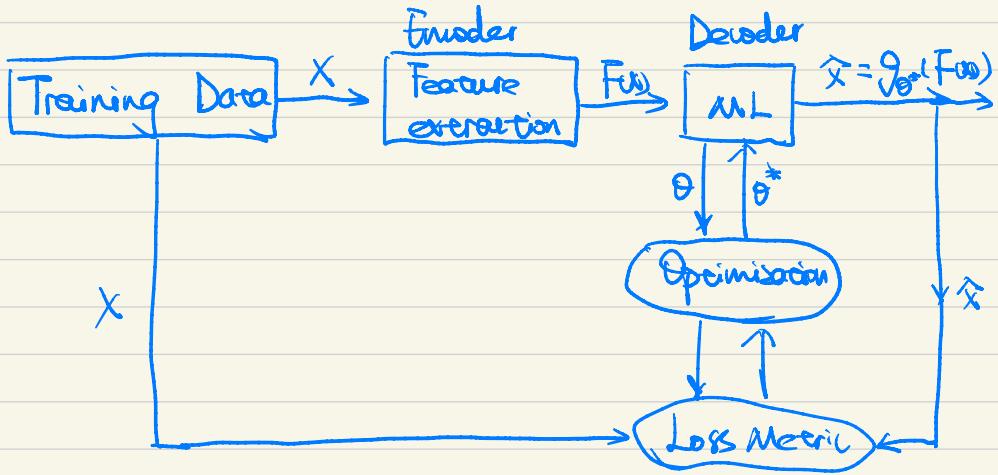
$$q_{ij} = \frac{\exp(-\|y_i - y_j\|_2^2)}{\sum_{k \neq i} \exp(-\|y_k - y_j\|_2^2)} \quad \text{SNE}$$

$$q_{ij} = \frac{1 + \|y_i - y_j\|_2^{-2}}{\sum_{k \neq i} 1 + \|y_k - y_j\|_2^{-2}} \quad t\text{-SNE} / \text{reduce crowding}$$

$$\min Q \text{KL}(P || Q) = \sum_{ij} P_{ij} \log \frac{P_{ij}}{q_{ij}}$$

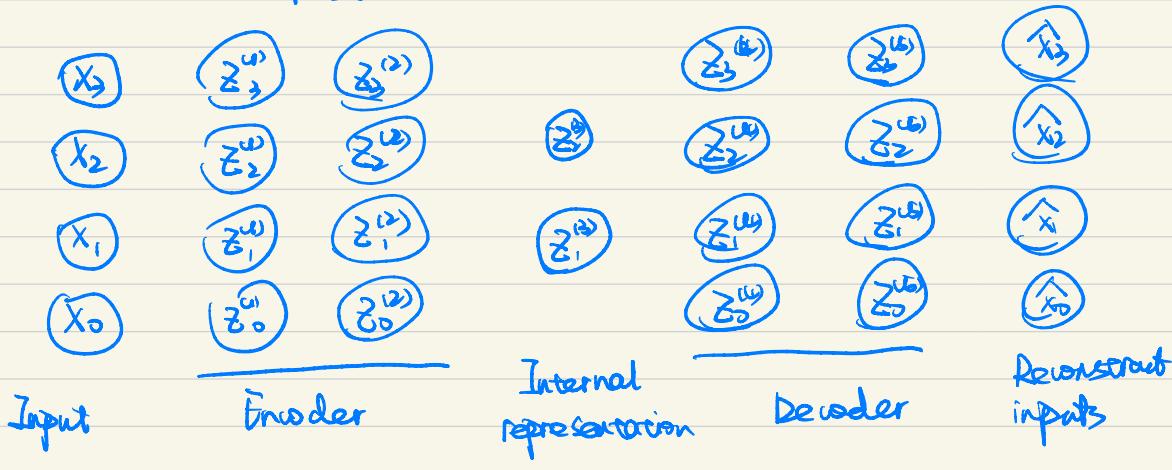
good \Rightarrow remain local pattern: KL not symmetric

Auto encoder



For PCA, keep both F and g linear

$$\begin{cases} F(X) = XU \text{ eigen} \\ \hat{X} = F(X)U^T \end{cases}$$



Min construction error \hat{X}, X

Support Vector Machine

S U M



Hard Margin SVM



idea: look for fat margin

few dichotomies better generalization

$$\text{Min } \frac{1}{2} w^T w$$

$$\text{s.t. } y_n (w^T x_n + b) \geq 1 \text{ for } n=1, \dots, N$$

Lagrangian:

$$L(w, b, \alpha) = \frac{1}{2} w^T w - \sum_{n=1}^N \alpha_n (y_n (w^T x_n + b) - 1)$$

$$\text{s.t. } \alpha_n \geq 0$$

$$\begin{array}{l} \text{plug } \frac{\partial L}{\partial w} = w - \sum_{n=1}^N \alpha_n y_n x_n = 0 \\ \frac{\partial L}{\partial b} = - \sum_{n=1}^N \alpha_n y_n = 0 \end{array}$$

$$L(w, b, \alpha) = L(\alpha) = \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N y_n y_m \alpha_n \alpha_m x_n^T x_m$$

$$\begin{array}{l} \text{max against } \alpha \\ \text{s.t. } \alpha_n \geq 0 \\ \sum \alpha_n y_n = 0 \end{array}$$

solve with a quadratic programming

$$\text{Min}_{\alpha} \frac{1}{2} \alpha^T J \alpha$$

$$J = \begin{bmatrix} y_1 y_1^T x_1 x_1 & \cdots & y_N y_N^T x_N x_N \\ \vdots & \ddots & \vdots \\ y_N y_1^T x_1 x_1 & \cdots & y_N y_N^T x_N x_N \end{bmatrix}$$

$$\text{s.t. } y^T \alpha = 0$$

$$0 \leq \alpha \leq \infty$$

replace $x_i^T x_j$
with $k(x_i, x_j)$
for kernel

$$\downarrow \alpha_n$$

$$w = \sum_{n=1}^N \alpha_n y_n x_n$$

↓ KKT condition

$\alpha_n > 0 \Rightarrow x_n$ is a support vector

Generalization depend on SV: $E[\text{Err}] = \frac{E[\# \text{ of SV}]}{N-1}$

Soft Margin SVM

$$\text{Min } \frac{1}{2} w^T w + C \sum \xi_m$$

$$\text{s.t. } y_n(w^T x_n + b) \geq 1 - \xi_m \\ \xi_m \geq 0$$

Lagrangian

$$L(w, b, \xi, \alpha, \beta) = \frac{1}{2} w^T w + C \sum \xi_m - \sum \alpha_n (y_n(w^T x_n + b) - 1 + \xi_m) - \sum \beta_n \xi_m$$

$$\frac{\partial L}{\partial w} = w - \sum \alpha_n y_n x_n = 0$$

$$\frac{\partial L}{\partial b} = - \sum \alpha_n y_n = 0$$

$$\frac{\partial L}{\partial \xi_m} = C - \alpha_m - \beta_m = 0$$

$$L(w) = \sum \alpha_n - \frac{1}{2} \sum \sum y_n y_m \alpha_n \alpha_m x_n^T x_m$$

$$\text{s.t. } 0 \leq \alpha_n \leq C$$

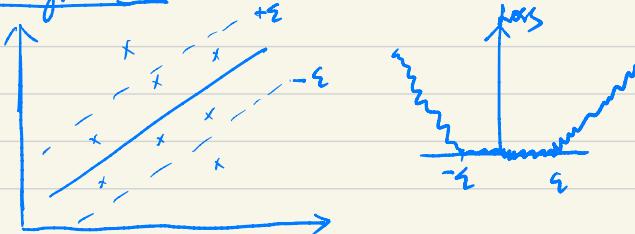
$$\sum \alpha_n y_n = 0$$

$$\Downarrow \alpha_n$$

$$w = \sum \alpha_n y_n x_n$$

$$\text{su } \left| \begin{array}{l} \xi_n = 0 \ (0 < \alpha_n < C) : y_n(w^T x_n + b) = 1 \\ \xi_n > 0 \ (\alpha_n = C) : y_n(w^T x_n + b) < 1 \end{array} \right.$$

SVM regression



Does not care error smaller than ϵ

Make function f as flat as possible \rightarrow less sensitive to noise

$$f(x) = w^T x + b$$

$$\text{Min } \frac{1}{2} w^T w + C \sum (\zeta_n + \zeta_n^*)$$

$$\text{s.t. } y_i - w^T x_i - b \leq \zeta_i + \zeta_{i*} \quad n=1 \dots N$$

$$w^T x_i + b - y_i \leq \zeta_i^* + \zeta_{i*}$$

$$\zeta_n, \zeta_n^* \geq 0$$

↓ Lagrangian

$$\begin{aligned} L := & \frac{1}{2} w^T w + C \sum (\zeta_n + \zeta_n^*) - \sum (\gamma_n \zeta_n + \gamma_n^* \zeta_n^*) \\ & - \sum \zeta_n (\epsilon + \zeta_n - y_n + w^T x_n + b) \\ & - \sum \zeta_n^* (\epsilon + \zeta_n^* + y_n - w^T x_n - b) \end{aligned}$$

where $\gamma_n, \gamma_n^* \geq 0$

$$\frac{\partial L}{\partial b} = \sum (\zeta_n - \zeta_n^*) = 0$$

$$\frac{\partial L}{\partial w} = w - \sum (\zeta_n - \zeta_n^*) x_n = 0$$

$$\frac{\partial L}{\partial \zeta_n} = C - \zeta_n - \gamma_n = 0$$

$$\text{dual: Max } L = -\frac{1}{2} \sum \sum (\zeta_n - \zeta_n^*) (\zeta_m - \zeta_m^*) x_n^T x_m - \epsilon \sum (\zeta_n + \zeta_n^*) + \sum y_n (\zeta_n - \zeta_n^*)$$

$$\text{s.t. } \sum (\zeta_n - \zeta_n^*) = 0$$

$$\downarrow \zeta_n$$

$$f(x) = \sum (\zeta_n - \zeta_n^*) x_n^T x + b$$

vanish for nonbinding

Ensemble

Learning



Ensemble: a group of predictors. Aggregate the prediction of a group of weak learners

1. Bagging: sample data and train model independently
(- variance) Usually homogeneous model

2. Boosting: learn sequentially in a adaptive way
(- bias) Usually homogeneous model

① AdaBoost

for t in $1, \dots, T$

a. weak learner: $h_t(x)$

$$h_t(x) = \arg \min_{\xi_t} = \arg \min \sum w_{it} \xi_{it}$$

b. update weight, learner

$$\downarrow \alpha_t = \frac{1}{2} \ln \left(\frac{1 - \xi_t}{\xi_t} \right)$$

$$F_t(x) = F_{t-1}(x) + \alpha_t h_t(x)$$

c. update error rate \rightarrow larger for failed sample

$$w_{i,t+1} = w_{i,t} e^{-y_i(\alpha_t h_t(x))}, \text{ renormalize}$$

② Gradient Boosting

fit new learner on previous residual

$$\uparrow f_{m+1}(x) = f_m(x) + h_m(x) = y_i$$

$$\downarrow h_m(x) = y_i - f_m(x)$$

3. Stacking
(- prediction error) learn a meta model to combine weak learners. Usually on heterogeneous weak learners

Tree Model



Decision Tree

Divide into subset

Measure of a node

classifier:

① Gini impurity

$$G_i = 1 - \sum_{k=1}^K P_{i,k}^2$$

proportion of class k in node i

② Entropy

$$H_i = - \sum_{\substack{k=1 \\ P_{i,k} \neq 0}}^K P_{i,k} \log(P_{i,k})$$

regression

$$MSE_h = \sum (y_j - \hat{y}_j)^2$$

CART Algorithm

Classification and regression tree - greedy method

At each step, choose feature and threshold to minimize

$$J = \frac{m_{left}}{m} G_{left} + \frac{m_{right}}{m} G_{right}$$

Random Forest

Decision tree + bagging

feature importance \rightarrow average depth

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 算的是 stochastic process To single realiz
- ② dependence 算的是 predict, 但也会有统计 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 → remove linear trend
2. Remove nonparametric regression on t
i.e.: kernel smoothing, polynomial, k-neighbor
→ remove trend / seasonality
3. Remove periodic regression on t → remove seasonality
4. Differentiation
1st → remove linear trend / expanding variance
2nd → remove polynomial trend / expanding variance
other interval → remove seasonality
...
5. Transformation
log → 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 ϕ_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

$$(\hat{\phi}_1, \hat{\phi}_2, \dots, \hat{\phi}_p, \hat{\sigma}_{w0}^2) = \tilde{\gamma}^{-1}(\hat{Y}_{X(0)}, \dots, \hat{Y}_{X(p)})$$

continuous mapping
Delta theorem
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

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

↓

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

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:

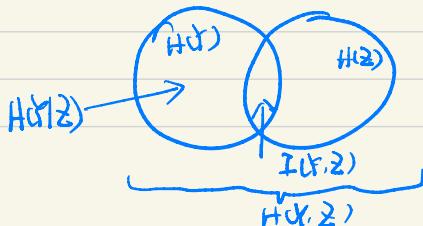
$$H(Y, Z) = - \sum_{y,z} P(y, z) \log P(y, z) \leq H(Y) + H(Z)$$

• Conditional Entropy:

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

• Mutual Information

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



Neural Network



Forward Propagation

$$\begin{cases} z^{[l]} = w^{[l]} a^{[l-1]} + b^{[l]} \\ a^{[l]} = g^{[l]}(z^{[l]}) \end{cases}$$

Vectorize over multiple samples

$$z^{[l]} = [z_1^{[l](1)} \dots z_1^{[l](m)}]$$

$b^{[l]}$ with broadcasting

$$\begin{cases} z^{[l]} = w^{[l]} A^{[l-1]} + b^{[l]} \\ A^{[l]} = g^{[l]}(z^{[l]}) \end{cases}$$

\dim

$w^{[l]} : (n^{[l]}, n^{[l-1]})$
 $b^{[l]} : (n^{[l]}, 1)$

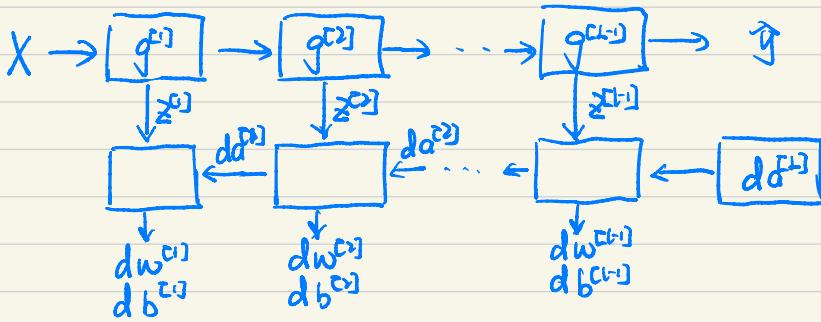
↓
Loop through l for forward propagation

Backward Propagation

$$\begin{cases} dz^{[l]} = da^{[l]} g^{[l]}'(z^{[l]}) \\ dw^{[l]} = dz^{[l]} a^{[l-1]T} \\ db^{[l]} = dz^{[l]} \\ da^{[l-1]} = w^{[l]T} dz^{[l]} \end{cases}$$

Vectorize

$$\begin{cases} dz^{[l]} = dA^{[l]} * g^{[l]}'(z^{[l]}) \\ dw^{[l]} = \frac{1}{m} dz^{[l]} A^{[l-1]T} \\ db^{[l]} = \frac{1}{m} np.sum(dz^{[l]}, axis=1, keepdims=True) \\ dA^{[l-1]} = w^{[l]T} dz^{[l]} \end{cases}$$



Exploding / Vanishing Gradients

$$\text{i.e. } g(z) = z \cdot b^L = 0$$

$$y = w^{[L]} w^{[L-1]} w^{[L-2]} \dots w^{[2]} w^{[1]} \underbrace{x}_{z^{[1]} \rightarrow z^{[L]}}$$

$\frac{z^{[1]} \rightarrow z^{[L]}}{z^{[2]} \rightarrow z^{[L]}}$

w large \rightarrow y explode \rightarrow gradient explode

w small \rightarrow y vanish \rightarrow gradient vanish

Solution:

Scale initialization

$$w^{[L]} = \text{np.random.randn(shape)} * \text{scale}$$

$$\text{scale} = \begin{cases} 0 & \sqrt{1/n^{[L-1]}} \\ \Theta & \sqrt{2/(n^{[L-1]} + n^{[L]})} \end{cases} \quad \text{Xavier}$$

Gradient Descent with momentum

almost always faster than gradient descent

$$\begin{cases} Vdw = \beta Vdw + u - \beta) dw \\ Vdb = \beta Vdb + u - \beta) db \end{cases}$$



RMSprop



On iteration t :

Compute dw, db on batch

$$\left\{ \begin{array}{l} Sdw = \beta Sdw + (1-\beta) dw^2 \\ Sdb = \beta Sdb + (1-\beta) db^2 \end{array} \right. \quad \begin{array}{l} \leftarrow \text{relative small} \\ \leftarrow \text{relative large} \end{array}$$

$$\left\{ \begin{array}{l} w = w - \alpha \frac{dw}{\sqrt{Sdw}} \\ b = b - \alpha \frac{db}{\sqrt{Sdb}} \end{array} \right.$$

Adam optimization adaptive moment estimation

- Initialization $Vdw = 0, Sdw = 0, Vdb = 0, Sdb = 0$
- On iteration t

$$\begin{matrix} \text{momentum} \\ + \end{matrix} \left\{ \begin{array}{l} Vdw = \beta_1 Vdw + (1-\beta_1) dw \\ Vdb = \beta_1 Vdb + (1-\beta_1) db \end{array} \right.$$

$$\begin{matrix} \text{RMSprop} \\ + \end{matrix} \left\{ \begin{array}{l} Sdw = \beta_2 Sdw + (1-\beta_2) dw^2 \\ Sdb = \beta_2 Sdb + (1-\beta_2) db^2 \end{array} \right.$$

$$w := w - \alpha \frac{Vdw}{\sqrt{Sdw} + \epsilon}, \quad b := b - \alpha \frac{Vdb}{\sqrt{Sdb} + \epsilon}$$

Batch Normalization

$$\text{normalize } z^{[t]} \text{ or sometime } \hat{z}^{[t]}$$

$$\hat{z}^{[t]} = (z^{[t]} - \mu) / \sqrt{\sigma^2 + \epsilon}$$

$$z^{[t]} = \gamma \hat{z}^{[t]} + \beta \quad \left\{ \begin{array}{l} \gamma = \sqrt{\sigma^2 + \epsilon} \\ \beta = \mu \end{array} \right.$$

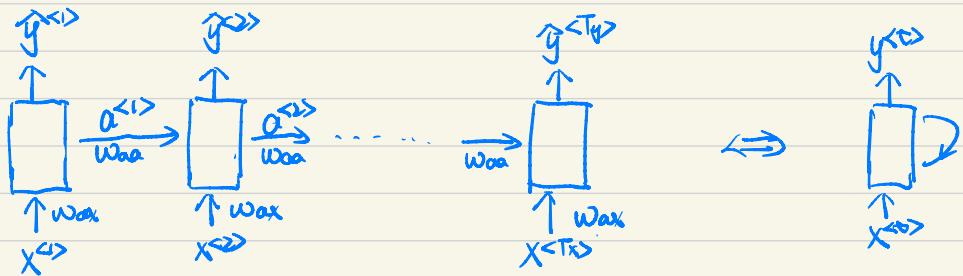
learnable parameter

RL

Reinforcement
Learning



Recurrent Neural Network



prediction based on new input $x^{<t>}$ as well as past information $a^{<t-1>}$

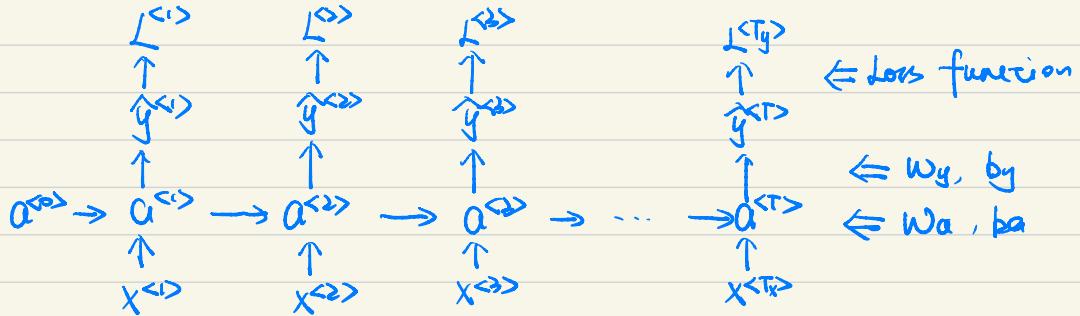
$$\left. \begin{array}{l} a^{<t>} = g(W_{aa}a^{<t-1>} + W_{ax}x^{<t>} + ba) = g(W_{aa}[a^{<t-1>} \cdot x^{<t>}] + ba) \\ y^{<t>} = g'(W_{ya}a^{<t>} + by) \\ a^{<0>} = \vec{0} \end{array} \right\}$$

$$W_a = [W_{aa} : W_{ax}]$$

$$[a^{<t-1>} \cdot x^{<t>}] = \begin{bmatrix} a^{<t-1>} \\ x^{<t>} \end{bmatrix}$$

parameter

Backpropagation

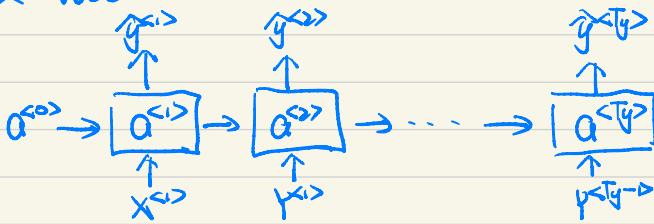


T_x and T_y can be different

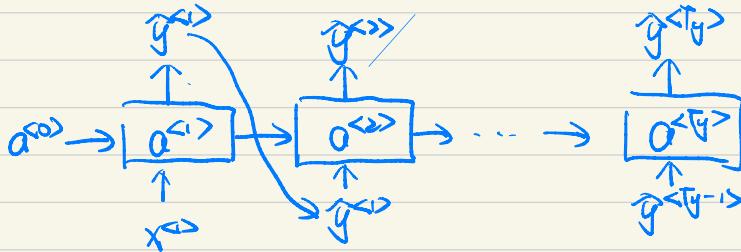
- ↳ many-to-many
- many-to-one
- one-to-many
- One-to-one

Language model

Training:

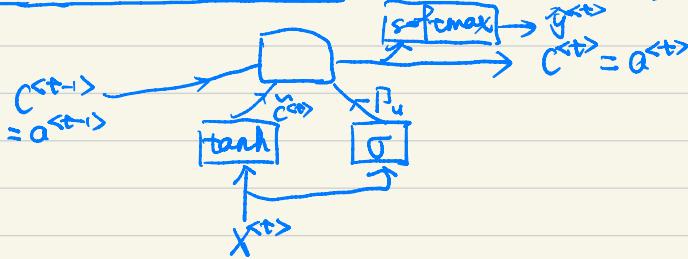


Sampling:



Gated recurrent unit (GRU)

→ better at capturing long term dependence



GRU

$$\begin{aligned} C^{<t>} &= \tanh(W_c[P_r \odot C^{<t-1>} + x^{<t>}] + b_c) \\ P_u &= \sigma(W_u[C^{<t-1>} + x^{<t>}] + b_u) \\ P_r &= \sigma(W_r[C^{<t-1>} + x^{<t>}] + b_r) \\ C^{<t>} &= P_u \odot C^{<t-1>} + (1 - P_u) \odot C^{<t-1>} \\ a^{<t>} &= C^{<t>} \end{aligned}$$

LSTM

$$\begin{aligned} C^{<t>} &= \tanh(W_c[a^{<t-1>} + x^{<t>}] + b_c) \\ P_u &= \sigma(W_u[a^{<t-1>} + x^{<t>}] + b_u) \quad \text{update} \\ P_f &= \sigma(W_f[a^{<t-1>} + x^{<t>}] + b_f) \quad \text{forget} \\ P_o &= \sigma(W_o[a^{<t-1>} + x^{<t>}] + b_o) \quad \text{output} \\ C^{<t>} &= P_u \odot C^{<t-1>} + P_f \odot C^{<t-1>} \\ a^{<t>} &= P_o \odot C^{<t>} \end{aligned}$$

Deep RNN

