

• Labels:

- Binary
- Multiclass
- Regression
- Structure prediction \rightarrow \square graph, sentence

• Protocol:

- Active
- Passive

• Data Assumption

- Batch \rightarrow typically iid
- Online \rightarrow iid / worst case / stochastic

* Why Theory

- Theory \heartsuit Practice

* Statistical Learning Theory (STL)

• Statistical

- x related to $y \rightarrow$ learn relationship \rightarrow distribution \mathcal{P} (population)
- we have access to samples of dist \rightarrow iid: all data from dist
independant sampling

\square Not good for spam filtering

\square Fair for hand writing recognition

$$- \mathcal{P} = \mathcal{P}_x \times \mathcal{P}_{y|x}$$

\swarrow marginal \searrow conditional

- $\hat{f}_n: X \mapsto Y$ mapping from dataset to function
 \hookrightarrow relation we wanted to learn

- ℓ : loss function $\ell: Y \times Y \mapsto \mathbb{R}$ measure of performance

- We're gonna average on all data, what my performance is

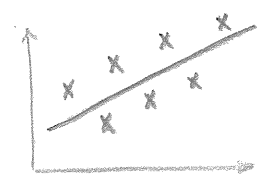
$$\mathbb{E}_{(x,y)} \{ \ell(\hat{f}_n(x), y) \} \rightarrow \text{Expectation of loss}$$

• Loss function

- Indicator Loss \rightarrow are they equal or not $\rightarrow \{0, 1\}$
- Square loss $\rightarrow (y - y')^2 \rightarrow$ regression cases
- Absolute loss $\rightarrow |y - y'|$

linear least squares

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n (y_i - \langle \beta, x_i \rangle)^2 \rightarrow \text{regression}$$



$$\hat{f}_n(x) = \langle \hat{\beta}, x \rangle$$

regularized least squares

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n (y_i - \langle \beta, x_i \rangle)^2 + \lambda \|\beta\|^2$$

So far we have nothing to do with problem (distribution)



- No best algorithms for all problems \rightarrow No free lunch

- Expected Loss: for any function $f: X \mapsto Y$

$$L(f) = \mathbb{E} e(f(x), y)$$

\mathcal{P} , unknown $\rightarrow L(f) = ?$

\mathcal{P} , known \rightarrow problem solved!

- Empirical loss: for any function

$$\hat{L}(f) = \frac{1}{n} \sum_{i=1}^n e(f(x_i), y_i)$$

\hat{L} : related to data

- what is Random here?

data are random i.i.d
draw from \mathcal{P}



Any function of
data is random
variable too



$\hat{L}(f) \rightarrow$ r.v.	
$\hat{f}_n \rightarrow$ r.v.	output of learning after seeing data
$L(\hat{f}_n) \rightarrow$ r.v.	
For a given $\mathcal{P}: X \mapsto Y$	
$L(f)$ is not not random!	

Gold Standard

- smallest expected loss
 \hookrightarrow not unique

$$f^* = \arg \min_f L(f)$$

\rightarrow Bayes Optimal

- Bayes Error $L(f^*) = \inf_f L(f)$

↳ \mathcal{P} unknown \rightarrow we cannot calculate

- We usually don't know the form of Bayes optimal function f^*
But two cases:

↳ Binary Classification $f^*(x) = \mathbb{I}\{\eta(x) \geq \frac{1}{2}\}$ $\eta(x) = \mathbb{E}[Y|X=x]$

↳ Regression $f^*(x) = \eta(x)$ $\eta(x) = \mathbb{E}[Y|X=x]$

\mathcal{P} : unknown \rightarrow value of f^* unknown but we know the form

• Start of STh \rightarrow Can we guarantee $L(\hat{f}_n) - L(f^*)$ to be small if n is large?

* Consistency

• Consistent: an algorithm that ensure

$$\lim_{n \rightarrow \infty} L(\hat{f}_n) = L(f^*)$$

• Good news: possible to achieve \rightarrow easy if X is finite / countable set
 \searrow not too hard if X is infinite and relation of X & Y is continuous

• There are always Bad news! \rightarrow we can't prove anything about $L(\hat{f}_n) - L(f^*)$
 \rightarrow need prior knowledge

No free Lunch Theorem

* No Free Lunch Theorem: if we don't have further assumptions

• for any alg $\hat{f}_n \rightarrow \mathcal{P}$ exist with $L(f^*) = 0$

if we don't give any assumption, we can't perform better than a coin tossing!

$$\mathbb{E} L(\hat{f}_n) \geq \frac{1}{2} - \epsilon$$

• for any alg $\hat{f}_n \rightarrow \mathcal{P}$ exist with $L(f^*) = 0$

we can't say how long we should wait for alg to perform good enough.

$$\lim_{n \rightarrow \infty} a_n = 0$$

Means \Rightarrow

$$\mathbb{E} L(\hat{f}_n) \geq a_n$$

• Is it really bad news?

- we always have domain knowledge

- Two ways of incorporating knowledge

↳ Direct knowledge → Modeling or Generative approach → we guess the distribution

↳ Indirect knowledge → Redefine the goal Statistics

$$L(\hat{f}_n) - \inf_{f \in F} L(f)$$

* Started by Vapnik → discriminative approach

* F encapsulates our inductive bias

↳ Cons / Pros → slide 33

→ we should know both → sometimes one of them is better

Linear Discriminant Analysis

$$\Sigma = I \rightarrow P(x|y=0) = (2\pi)^{-k/2} \exp\left\{-\frac{1}{2}\|x-\mu_0\|^2\right\} \Rightarrow \mu_0, \mu_1: \text{unknown}$$

$$P(x|y=1) = (2\pi)^{-k/2} \exp\left\{-\frac{1}{2}\|x-\mu_1\|^2\right\}$$

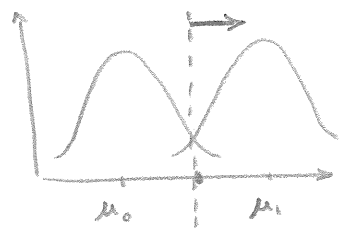
Best Classifier $f^* = I_{P(y=1|x) \geq 1/2}$ → linear : $P(x|y=1) \geq P(x|y=0)$

↳ Bayes

$\Sigma_1 \neq \Sigma_2$ → Quadratic Discr. Anal.

$$I\{q(x) \geq 0\}$$

$q(x)$ ~ quadratic function



* Bias - Variance Trade off

• Choosing a model (assumption) → Bias

• How do we choose F?

$$L(\hat{f}_n) - L(f^*) = \underbrace{L(\hat{f}_n) - \inf_{f \in F} L(f)}_{\text{Est. Err}} + \underbrace{\inf_{f \in F} L(f) - L(f^*)}_{\text{Approx. Err}}$$

- F larger → smaller approx err

- n ↑ → smaller estimation err. & no effect on approximation err.

• if we guessed f^* so $f^* \in F \rightarrow L(\hat{f}_n) - L(f^*) = L(\hat{f}_n) - \inf_{f \in F} L(f)$

* Occam's Razor

- Principle for choosing the simplest theory or explanation
- Help to try understanding the complexity
- We look for alg \hat{f}_n that with n samples have small $L(\hat{f}_n) - \inf_{f \in F} L(f)$
 - ↳ smaller than desired accuracy
 - ↳ with fixed set of $F \rightarrow$ rich enough

• Bounds

- In Expectation $E \{ L(\hat{f}_n) - \inf_{f \in F} L(f) \}$

- In Probability $\mathbb{P} \{ L(\hat{f}_n) - \inf_{f \in F} L(f) \geq \epsilon \} < \psi(n, \epsilon)$
 ↳ control the tail

or $\mathbb{P} (L(\hat{f}_n) - \inf_{f \in F} L(f) \geq \Phi(\delta, n)) < \delta$

$\delta := \psi(n, \epsilon) \rightarrow$ solving for ϵ

$\Phi(\delta, n) \rightarrow$ we want in this lecture \rightarrow "the rate"

- Sample Complexity \rightarrow what is the # of examples algorithm need to reach ϵ accuracy with $(1-\delta)$ probability

$\mathbb{P} (L(\hat{f}_n) - \inf_{f \in F} L(f) \geq \Phi(\delta, n)) < \delta$

$\epsilon := \Phi(\delta, n) \rightarrow$ solve for n

$\mathbb{P} (L(\hat{f}_n) - \inf_{f \in F} L(f) \geq \epsilon) < \delta \quad n \geq n(\epsilon, \delta)$

$n \downarrow \Rightarrow$ faster rate to reach ϵ -guarantee

$O(\frac{1}{\sqrt{n}}) \Rightarrow O(\frac{1}{\epsilon^2}), O(\frac{1}{n}) \Rightarrow O(\frac{1}{\epsilon})$

* Clarify

(A)

(D) Look at function and difference to empiric data \rightarrow which one is better

slide 44

(B) To guarantee the performance of f

(E) Upper bound for (B) Model Selection \rightarrow grow F

(C)

* Bounds

- $L(\hat{f}_n) - \hat{L}(\hat{f}_n) \leq \sup_{f \in F} \{L(f) - \hat{L}(f)\} \Rightarrow (D) \text{ implies a bound on } (C)$

- Bound on (B)? \rightarrow (C) or (D)? depends on what the alg does

- ERM: Empirical Risk Minimization $\hat{f}_n = \arg \min_{f \in F} \hat{L}(f)$

- Regularized ERM $\hat{f}_n = \arg \min_{f \in F} \hat{L}(f) + \text{Pen}_n(f)$

- $f_F = \inf_{f \in F} L(f) \rightarrow$ one of the best choices for f

$$L(\hat{f}_n) - L(f_F) \leq \underbrace{\{L(\hat{f}_n) - \hat{L}(\hat{f}_n)\}}_{(C)} + \underbrace{\{\hat{L}(\hat{f}_n) - \hat{L}(f_F)\}}_{\leq 0} + \underbrace{\{\hat{L}(f_F) - L(f_F)\}}_{(*)}$$

$$\leq (C) + (*)$$

• Hoeffding Inequality

$w_i \rightarrow$ i.i.d
 $P(a \leq w \leq b) = 1$

$$\Rightarrow P(\mathbb{E}w - \frac{1}{n} \sum_{i=1}^n w_i > \epsilon) \leq \exp\left(-\frac{2n\epsilon^2}{(b-a)^2}\right)$$

gaussian tails are enforced by $\epsilon \Rightarrow \epsilon \uparrow \rightarrow$ more gaussian tails

using central limit theorem

$$P\left(\frac{1}{n} \sum_{i=1}^n w_i - \mathbb{E}w > \epsilon\right) \leq \exp\left(-\frac{2n\epsilon^2}{(b-a)^2}\right)$$

$$\Rightarrow P(|L(f_F) - \hat{L}(f_F)| > \epsilon) \leq 2 \exp\left(-\frac{2n\epsilon^2}{(b-a)^2}\right)$$

$$a \leq \ell(f_F(x), y) \leq b$$

- $P(\mathbb{E}L(\hat{f}_n(x), y) - \underbrace{\frac{1}{n} \sum \ell(\hat{f}_n(x_i), y_i)}_{\text{empirical}} > \epsilon) \leq 2 \exp(\sim)$

$\ell(\hat{f}_n(x), y) \rightarrow$ memory, test on same training set \rightarrow cannot generalize
 \rightarrow biased estimate of $\mathbb{E} \ell(\hat{f}_n(x), y) \rightarrow$ overfitting

out of sample performance $\left\{ \begin{array}{l} \mathbb{E}w_o = \mathbb{E} \ell(\hat{f}_n(x), y) \\ \mathbb{E}w_i = \mathbb{E} \ell(\hat{f}_n(x_i), y_i) \end{array} \right\} \neq \Rightarrow \text{gap} \rightarrow \text{overfitting}$

- ERM $\hat{f}_n \rightarrow$ memory, perfect recall of data \rightarrow cannot generalize
- $\sup_{f \in F} \{ \underbrace{L(f)}_1 - \underbrace{\hat{L}(f)}_{0 = \text{perfect fit}} \} = 1 \rightarrow$ (D): called uniform deviation
- cannot generalize

* Uniform Deviation

$$\sup_{f \in F} \left\{ \mathbb{E}_{x,y} \ell(f(x), y) - \frac{1}{n} \sum \ell(f(x_i), y_i) \right\}$$

$F = \{f_0\} \rightarrow$ have good bound, always answer the same
 $F = \text{rich} \rightarrow$ maybe has bad bound

$z_i = (x_i, y_i) \rightarrow$ training pair

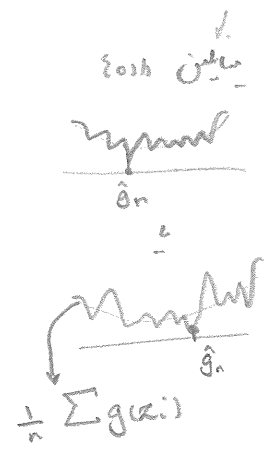
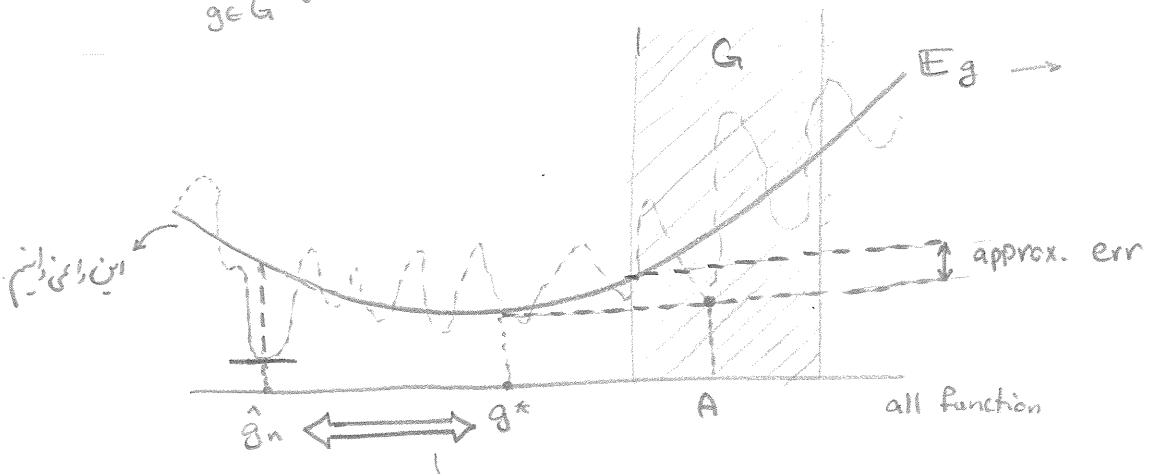
$$g(z) = \ell(f(x), y)$$

$$\hat{g}_n = \ell(\hat{f}_n(\cdot), \cdot)$$

$$g_G = \ell(f_F(\cdot), \cdot)$$

$$g^* = \underset{g}{\operatorname{argmin}} \mathbb{E} g(z) = \ell(f^*(\cdot), \cdot) \rightarrow \begin{cases} \text{Bayes optimal loss function} \\ \text{best overall function} \end{cases}$$

$$\star = \sup_{g \in G} \left\{ \mathbb{E} g(z) - \frac{1}{n} \sum g(z_i) \right\}$$



این را می دانیم
 بهترین نتیجه درست آمده از داده؟
 نتیجه به سبب آماری ممکن است
 overfitting

\hat{g}_n به تنهایی "A" مستقل می شود (در نتیجه درست زده شده) \rightarrow می خواهیم به اندازه

در پیچیده G را کوچک کنیم
 کافی G را کوچک کنیم.

• Stochastic Process: collection of random var. index by some set.

- Empirical Process \leadsto stochastic process by function G

- Uniform Law of Large Number

$$\sup_{g \in G} \left| \mathbb{E}g - \frac{1}{n} \sum g(z_i) \right| \xrightarrow{n \rightarrow \infty} \infty$$

- question: How big can "G" be?

- reminder: Boole's Inequality $\mathbb{P}(\cup A_j) \leq \sum \mathbb{P}(A_j)$

$$\mathbb{P}(\exists g \in G : \mathbb{E}g - \frac{1}{n} \sum g(z_i) > \epsilon) \leq \sum_j \mathbb{P}(\mathbb{E}g_j - \frac{1}{n} \sum g_j(z_i) > \epsilon)$$

$$\left\{ \begin{array}{l} G = \{g_1, g_2, \dots, g_N\} \quad \# \text{ of function} \\ \mathbb{P}(a \leq g(z_i) \leq b) = 1 \end{array} \right.$$

$\mathbb{P}(\dots)$

- How many example do we need?

$$\mathbb{P}\left(\sup_{g \in G} \left\{ \mathbb{E}g - \frac{1}{n} \sum g(z_i) \right\} > (b-a) \sqrt{\frac{\log N + \log(1/\delta)}{2n}}\right) \leq \delta$$

--- slide 56

• Message

$\left\{ \begin{array}{l} \rightarrow \text{with prob at least } 1-\delta \text{ for any } g \in G \Rightarrow \mathbb{E}g - \frac{1}{n} \sum g(z_i) \leq \epsilon \\ \rightarrow \text{for any } g \in G \text{ with prob at least } 1-\delta \Rightarrow \mathbb{E}g - \frac{1}{n} \sum g(z_i) \leq \epsilon \end{array} \right.$

• Weighted Union Bound: Countable Class

$$\sum_{g \in G} \omega(g) \leq 1$$

$$L(f) - \hat{L}(f) \leq (b-a) \underbrace{\sqrt{\frac{\log \frac{1}{\omega(f)} + \log \frac{1}{\delta}}{2n}}}_{pen_n(f)}$$

$$L(\hat{f}_n) - L(f_F) \leq \{L(\hat{f}_n) - \tilde{L}(f_n) - \text{pen}_n(f_n)\} + \{\hat{L}(\hat{f}_n) + \text{pen}_n(\hat{f}_n) - \hat{L}(f_F) - \text{pen}_n(f_F)\} \leq 0$$

$$+ \{\hat{L}(f_F) - L(f_F)\} + \text{pen}_n(f_F)$$

$$\leq \sup_{f \in F} \{L(f) - \hat{L}(f) - \text{pen}_n(f)\} + \{\hat{L}(f_F) - L(f_F)\} + \text{pen}_n(f_F)$$

$$L(\hat{f}_n) - L(f_F) \leq \{\hat{L}(f_F) - L(f_F)\} + \text{pen}_n(f_F)$$

$$\leq 2(b-a) \sqrt{\frac{\log \frac{1}{\omega(f_F)} + \log \frac{1}{\delta}}{2n}}$$

• Uncountable Class: Compression Bound

algorithm \hat{f}_n depend on dataset $S \rightarrow \hat{f}_n = \hat{f}_n[S]$

C_k select a subset of useful examples $\rightarrow \hat{f}_n[S], \hat{f}_k[C_k(S)] \rightarrow$ not general

$$L(\hat{f}_n) - \hat{L}(\hat{f}_n) = \mathbb{E} \ell(\hat{f}_k[C_k(S)](x_i, y_i)) - \frac{1}{n} \sum \hat{f}_n[C_k(S)](x_i, y_i)$$

$$\leq \max \quad \text{slide 6L}$$

most of data are now out of samples

$$L(\hat{f}_n) - \inf_{f \in F} L(f) \leq 2(b-a) \sqrt{\frac{k \log(\frac{en}{k}) + \log \frac{1}{\delta}}{2n}} + 2 \frac{(b-a)k}{n}$$

$F_{k,n}$ has cardinality at most $\binom{n}{k} \leq (\frac{en}{k})^k$

* Stability

• change of a single example don't change the algorithm too much

$$\hat{f}_n = \arg \min \frac{1}{n} \sum (f(x_i) - y_i)^2 + \lambda \|f\|_k^2$$

↪ stable \rightarrow Reproducing kernel Hilbert space

* Uniform Convergence

$z_1, \dots, z_n \rightarrow$ iid from \mathcal{P}

$\epsilon_1, \dots, \epsilon_n \rightarrow$ iid with $P(\epsilon_i = -1) = P(\epsilon_i = +1) = 1/2 \rightarrow$ coin flip
Rademacher r.v.

$$\mathbb{E} \sup_{g \in G} \left\{ \mathbb{E} g(z) - \frac{1}{n} \sum g(z_i) \right\} = \mathbb{E} \sup_{z_{1:n}} \left\{ \frac{1}{n} \sum g(z_i) - \frac{1}{n} \sum g(\epsilon_i z_i) \right\}$$

comparison of empirical with model \Rightarrow comparison of empirical 1 to empirical 2

* Symmetrization

$$\mathbb{E} L(f_n) - \inf_{f \in F} L(f) \leq \mathbb{E} \sup_{f \in F} L(f) - \hat{L}(f) \leq 2 \underbrace{\mathbb{E} \mathbb{E}_{\text{data } \epsilon_{1:n}} \sup_{f \in F} \frac{1}{n} \sum \epsilon_i g(z_i)}_{\text{Rademacher Complexity Bound}}$$

data points z_i : $\left\{ \begin{array}{l} \text{function 1 : } f(z_1), \dots, f(z_n) \\ \text{function 2 : } f'(z_1), \dots, f'(z_n) \\ \dots \end{array} \right\}$ $F|_{x_{1:n}} \subseteq \{0,1\}^n$
 ↪ subset of nD space

$$v = \begin{bmatrix} f(z_1) \\ \vdots \\ f(z_n) \end{bmatrix} \Rightarrow v \cdot \epsilon = \sum_i \epsilon_i \cdot v_i$$

VC classes

Half spaces → half predict 0, half predict 1

* Large Margin Theory for Classification

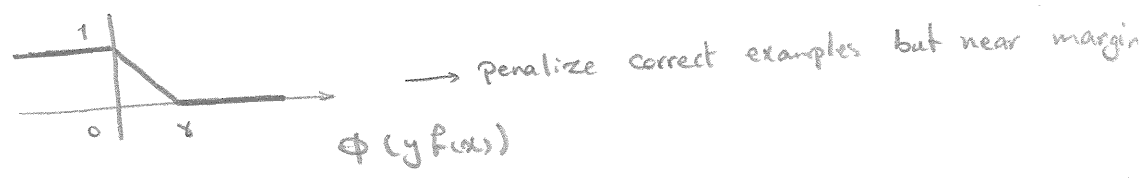
- $I(F) = \{I_{\{f \geq 0\}} : f \in F\}$
- $f(x) = f_w(x) = \langle w, \psi(x) \rangle \rightarrow \psi(x) = \text{kernel} \rightarrow$ VC dimension is high
 ↘ but SVM working on this basis work well

• Hard margin → $\exists f \in F : \forall i, y_i f(x_i) \geq \gamma$

$\exists f \in F : \frac{\text{misclassification}}{n} \rightarrow$ hope to be small



• Soft margin



Surrogate loss

$$\mathbb{E} \sup_{f \in F} \{L_\phi(f) - \hat{L}_\phi(f)\}$$

if ϕ is L -lipschitz → $R_n(\phi(F))$ is bounded

* Properties of Rademacher Complexity

- 1.
2. Convex hull of function → don't increase R
3. $\hat{R}_n(CF) = |\text{cl}| \hat{R}(F)$
- 4.

- Easier to obtain sample complexity for Learning Tree, NN, ...
 - ↳ using Rademacher complexity

* Sequential Prediction

• Motivation

- Future look like the past
- Get rid of iid

□ www.secd.ucsd.edu/~mindreader → human cannot make a random sequence we are finite state machines

↳ input: $y_1, y_2, \dots \in \{0,1\}$ → one by one

↳ output: x_1, x_2, \dots → prediction

↳ goal → make average of correct predictions large

↘ proportion of correction prediction is larger than proportion of 1's or 0's
 ↘ make less mistakes than two predictors → { one always says 1
 experts { one always says 0
 strategists

□ Spam Detection

x_1, \dots, x_n → emails, revealed one-by-one

↳ not iid

↳ spammers adapt to what you do → adversarial

• Supervised Learning

- X input, Y output

- Regret = you are close to some hypothesis (# of incorrect predictions)
 - if you use that f you suffer that much

- Goal: keep regret small for any sequence of (x_i, y_i)

• Online Convex and Linear Optimization

- loss function is convex

□ squared loss $\ell(f, (x, y)) = (\langle f, x \rangle - y)^2$

hinge loss $\ell(f, (x, y)) = \max\{0, 1 - y \langle f, x \rangle\}$

- $\text{Reg}_n = \frac{1}{n} \sum_t \ell(f_t, z_t) - \inf_{f \in F} \frac{1}{n} \sum_t \ell(f, z_t)$

- Gradient Descent: $f_{t+1} = f_t - \eta \nabla \ell(f_t, z_t)$

- Alg. guarantee $\text{Reg}_n \leq O\left(\frac{1}{\sqrt{n}}\right)$

- Function strongly convex → grow faster than quadratic function

- $\forall f^* \in F : \text{Reg}_n \leq O\left(\frac{\log n}{n}\right)$

- Use for iid :

$\text{Reg}_n \leq R_n$

$\mathbb{E} L(\bar{f}) - \inf_{f \in F} L(f) \leq R_n$

Distribution \mathcal{P} ~ unknown

↳ Realize in "Pegasos" alg.

- Pegasos (= SVM in linear case)

$\hat{f}_n = \arg \min_{f \in \mathbb{R}^d} \frac{1}{m} \sum \max\{0, 1 - y_i \langle f, x_i \rangle\} + \frac{\lambda}{2} \|f\|^2$

$l(f, z) = \max\{0, 1 - y \langle f, x \rangle\} + \frac{\lambda}{2} \|f\|^2 \rightarrow$ hinge loss func.

goal: $\min_f \mathbb{E} l(f, z)$

guarantee: $\mathbb{E} l(\bar{f}, z) - \inf_f \mathbb{E} l(f, z) \leq O\left(\frac{\log n}{n}\right)$

alg: for $t = 1, \dots, n$

Choose a random example (x_t, y_t)

set $\eta_t = \frac{1}{\lambda t}$

if $y_t \langle f_t, x_t \rangle < 1$

$f_{t+1} = (1 - \eta_t \lambda) f_t + \eta_t x_t y_t$

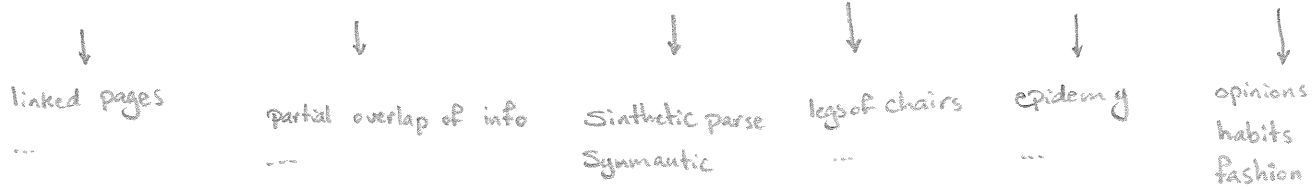
else

$f_{t+1} = (1 - \eta_t \lambda) f_t$

* Motivation

- iid assumption \rightarrow one type of object, no relation
- Real application \rightarrow multiple objects, heavily related

☑ web search, information extraction, NLP, perception, Medical Diagnosis, Social Net.



* Costs & Benefits

- Benefits \rightarrow Better predictive Accuracy
- Benefits \rightarrow Better understanding of Domains
- Benefits \rightarrow Growth path of ML
- Costs \rightarrow Learning is much harder
- Costs \rightarrow iid \rightarrow 100 features \rightarrow vector of 100
- Costs \rightarrow SRL \rightarrow 100 features \rightarrow things has interaction
- Costs \rightarrow Inference becomes a crucial issue
- Costs \rightarrow \hookrightarrow even harder problem than learning
- Costs \rightarrow Greater complexity for user

* Goal

- goal: learn from non-iid data as easy as iid ones
- progress \rightarrow "close enough" to goal
- progress \rightarrow burgeoning research area
- progress \rightarrow Open Source, easy to use
- old & new questions

* Plan

- elements: probability, logic, learning and inference
- put them together
- lots of application

* Probabilistic Inference : Markov Networks

- Undirected graphical model → structure of dependencies
- Potential functions defined over cliques → Markov Blanket

$$P(x) = \frac{1}{Z} \prod_c \phi_c(x_c)$$

$$\phi_c \geq 0$$

↳ no restriction $\sum \phi \neq 1$

$$Z = \sum_c \prod_c \phi_c(x_c)$$

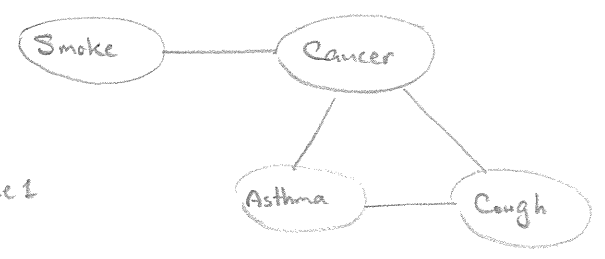
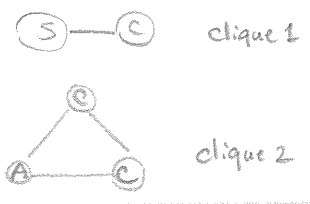
→ partition function : the most important part of Markov Net.

- Conditional independence advantage
- Size of Tables are exponential to size of cliques
- Log linear model

$$P(x) = \frac{1}{Z} \exp \left(\sum_i w_i f_i(x_i) \right)$$

✓ $f_1(\text{Smoking}, \text{Cancer}) = \begin{cases} 0 \\ 1 \end{cases}$
 $w_1 = 1.5$

Sm	Can	Φ
0	0	4.5
0	1	4.5
1	0	2.7
1	1	4.5



• Hammersley - Clifford Theorem

- Distribution is strictly positive $P(x) > 0$
 Graph Encodes conditional potentials } ⇒ Dist = \prod potentials over cliques of graph

- Inverse → Markov Networks = Gibbs distribution is True

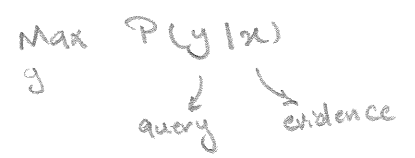
• MN vs. Bayesian Networks → Slide 12

- ↳ Real world have cycles → MN ✓
- ↳ Z is unknown → MN ☹️
- ↳ Indep check → MN is easier ✓ → easier to interpret

• MAP / MPE Inference

↳ Most Probable Explanation

- Goal : Find most likely state of world given evidence

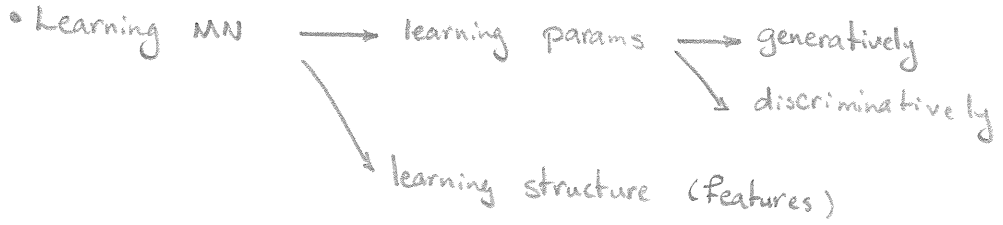


- Iterated conditional modes → ⚡ Fast
- ☹️ sensitive to local optima

- Simulate annealing → very slow
- for large iteration, it find the global min

- Graph Cuts
- Belief Propagation (max-product)

* Statistical Learning



• Generative weight learning

- Convex problem → no local maxima
- numerical optimization → gradient or 2nd order

$$\frac{\partial}{\partial \omega_i} \log P_w(x) = \underbrace{n_i(x)}_{\text{\# of times feature } i \text{ is true in data}} - \underbrace{E_w[n_i(x)]}_{\text{expected \# of times feature } i \text{ is true according to model}}$$

underestimate: model < real → $\frac{\partial}{\partial \omega_i} \geq 0 \rightarrow f_i \downarrow$

- requires inference at each step → too slow!

• Pseudo likelihood → if the problem you have is too hard, solve an easy one!

- Fake likelihood → easy to compute

$$PL(x) = \prod_i P(x_i | \text{neighbours}(x_i))$$

→ Markov Blanket

- no need to inference each step
- ✓ hold true for many cases: machine vision
- may not work well for long inference chains

• Discriminative weight learning

- Conditional random fields
- Maximize conditional likelihood of query (y) given evidence (x)

$$\frac{\partial}{\partial \omega_i} \log P_w(y|x) = n_i(x, y) - \underbrace{E_w[n_i(x, y)]}$$

can be approximated by counts in MAP state of y given x

• Other weight learning

- Generative: Iterative scaling
- Discriminative: Max margin

• Structure learning

- Start with atomic features
- Join features \rightarrow Improved score? \rightarrow keep it!
- Problem: calculate all weights for every new feature
- Approx: keep the weights of previous features constant

* Logical Inference

• First-Order Logic

- Constants, variables, functions, predicates
- Literal \rightarrow predicate
- Clause \rightarrow Disjunction of literals
- Grounding
- World

• Inference in First-order Logic

- Propositionalization
- Model checking \rightarrow Backtracking \rightarrow DPLL
- Stochastic Local Search \rightarrow Walk SAT

• Satisfiability

- Inputs \rightarrow sets of clauses \rightarrow CNF
- Output \rightarrow Truth assignment
- Solution \rightarrow Search
 - \rightarrow NP complete
 - \rightarrow Most SAT problems are easy
 - \rightarrow Hard region: narrow range of $\frac{\# \text{ clauses}}{\# \text{ variable}}$

• Backtracking

- Assign truth val. by DFS \rightarrow Dead end? \rightarrow change last assignment
- Empty set of clauses \rightarrow Success
- Empty clause \rightarrow Failure
- Assigning a variable \rightarrow delete false literals
 - \rightarrow satisfies clauses
- Improvements \rightarrow Unit propagation \rightarrow Fast
 - \rightarrow Pure literals

• Stochastic local Search

- Completely assign → random initial → flip vars in unsatisfied clauses
- minimize # unsatisfied (Hill climbing) → random flip (avoid local optima)
- multiple restarts

✓ Walk SAT → Faster DPLL
 ↘ if cannot find result it doesn't mean there's not one!

* Inductive Logic Programming

• Rule Induction

- Examples → x : attributes — Algorithm
 ↘ y : concept
- Goal → coverage of + samples
 ↘ not covering of any - sample
- Eval → Accuracy
 ↘ Coverage

✓ Spam filtering with bag-of-words model

• First - Order Rule Induction

- literals to add → predicates or their negations
 ↘ must include at least one var.
 ↘ changing the # of grounding of rule
- Eval

* Putting things together

✓ KB model construction

- Stochastic logic Programs
- Probabilistic Relational Models

• keys

- Logical Language → First order, Horn Clauses
- Probabilistic Language → Bayes Net, Markov Net, PCFG
- Types of Learning → Generative / Discr., Structure / Params, knowledge rich / poor
- Type of inference

• KB Model Construction

- Logical → Horn

- Prob → Bayes Net

- Learning → ILP + EM

- Inference → Partial Grounding
↳ Belief Prop

Ground Atom: Node

Head of Clause: Child Node

Body of Node: Parent Node

Combination function: Noisy OR

> 1 clause with same head

• Stochastic Logic Programs

- Logical → Horn

- Prob → Probabilistic CFG

- Learning → ILP

↳ "Failure Adjusted" EM

- Inference → Do all proof

↳ Add probabilities

Attach probs. to clauses

$\sum (\text{clauses with same head}) = 1$

• Probabilistic Relational Model

- Logical → Frame

- Prob → Bayes Net

BN temple for each class of object

Obj 1 attr. can be dependant on Obj 2 attr.

Binary relations

relations are independent

- Learning → Params: Closed form (EM for missing data)

↳ Structure: "Tiered" BN

- Inference → Full Grounding

↳ BP

• Relational Markov Nets

- Logical → SQL query

- Prob → Markov Net

SQL query define clique

Potential for each query

No uncertainty over relations

- Learning → Discriminative weight Learning

↳ No structure learning

- Inference → Full Grounding

↳ Belief Propagation

• Bayesian Logic

- Logic \rightarrow First order
- Prob \rightarrow BN \rightarrow BLOG specifies how to generate relational world

allows objects to be uncertain

• Markov Logic

- Logic \rightarrow First order
- Prob \rightarrow MN \rightarrow Syntax: First order formula + weights
 \searrow Semantics: Template for Markov net features
- Learning \rightarrow Params: Gen / Discr
 \searrow Structure: ILP with arbitrary clauses and MAP scores
- Inference \rightarrow MAP: weighted satisfiability
 \searrow Marginal: MCMC \rightarrow moves proposed by SAT solver
 \searrow Partial grounding + Lazy inference

- Logical KB: a simple false clause in world \rightarrow make it impossible (hard)

Markov Logic: a false clause make world less probable (soft)

- in Markov Logic each formula has a weight of truth

☑ Smoking cause cancer \rightarrow Logic: $\forall x \text{ Smoke}(x) \Rightarrow \text{Cancer}(x)$
 \searrow Markov: 1.5 $\forall x \dots$

- $R(x,y)$ results to $\begin{cases} R(x,y) \\ R(y,x) \\ R(x,x) \\ R(y,y) \end{cases}$

☑ Friend(A,B) \rightarrow Friendship is asymmetric

- Cliques are shown by relations and atoms

- Statistical Model: Special case of Markov Logic with zero-arity

- First order Logic: " " " " " " infinite weight

- MAP Inference

$$\underset{\theta}{\operatorname{argmax}} \frac{1}{Z_{\theta}} \exp\left(\sum_i w_i n_i(\alpha, y)\right) \Rightarrow \underset{y}{\operatorname{argmax}} \sum_i w_i n_i(\alpha, y) \Rightarrow \text{weighted maxSAT problem}$$

☑ MaxWalkSAT

↳ Problem: Memory explosion → Solution: ground clauses lazily → sparseness

☑ Lazy SAT alg.

- Computing prob

↳ $P(\text{Formula} \mid \text{MLN}, C) = ? \rightarrow$ MCMC sample & check formula

↳ $P(f_1 \mid f_2; \text{MLN}, C) = ? \rightarrow$ make world \rightarrow satisfy f_2 ? \rightarrow satisfy f_1 ? yes ✓

↙ use lifted inference \rightarrow active field

↘ construct min subset of net to satisfy $f_2 \rightarrow$ MCMC \rightarrow ...

↳ Problem \rightarrow MCMC fails most of the time \rightarrow deterministic dependency

↳ MCMC slow \rightarrow near deterministic dependency

↳ Solution \rightarrow Combine MCMC + walk SAT

☑ MC-SAT algorithm

* Applications of SRL

• Running Alchemy

- Programs

- Options

- MLN File

- Database File

↳ Robust Programs
Smaller Codes

☑ Uniform Dist: Flipping Coin

empty MLN \rightarrow Type: flip = {1, ..., 20}

Predicate: Heads (flip)

$$P(\text{Heads}(f)) = \frac{\frac{1}{2} e^{0}}{\frac{1}{2} e^0 + \frac{1}{2} e^0} = \frac{1}{2}$$

number of features = 0

☑ Binomial Dist: Unfair coin

Type: flip = {1, ..., 20}

Pred.: Heads (flip)

using unit clause

Formula: Heads (f)

Weight: log odds of head $\rightarrow w = \log\left(\frac{p}{1-p}\right)$

$$P(\text{Heads}(f)) = \frac{\frac{1}{2} e^w}{\frac{1}{2} e^w + \frac{1}{2} e^0} = \frac{1}{1 + e^{-w}} = p$$

☑ Multi Nominal Dist: Dice

Type: throw = {1, ..., 20}
face = {1, ..., 6}

Predicate: outcome {throw, face!}

Predicate: outcome (throw, face)

Formula: Outcome (t, f) $\wedge f \neq f' \Rightarrow$!outcome (t, f'), ...

↑ hard, need shortcut, use blocking

Formula: $\left. \begin{matrix} \text{outcome } (t_1, \dots) \\ \text{outcome } (t_2, \dots) \\ \text{outcome } (t_3, \dots) \\ \dots \end{matrix} \right\} \Rightarrow \text{outcome } (t, f)$

Logistic Regression

$$\log \left(\frac{P(C=1 | F=f)}{P(C=0 | F=f)} \right) = a + \sum b_i f_i$$

Type: $\text{obj} = \{1, \dots, n\}$

Query Pred: $C(\text{obj})$

Evidence Pred: $F_i(\text{obj})$

Formulas: $a \ C(x) \rightarrow a \text{ is weight}$
 $b_i \ F_i(x) \wedge C(x) \rightarrow b \text{ is weight}$

Resulting Dist: $P(C=c | F=f) = \frac{1}{Z} \exp(a_c + \sum b_i f_i \cdot c)$

$$\log \left(\frac{P(C=1 | F=f)}{P(C=0 | F=f)} \right) = \log \left(\frac{\exp(a_1 + \sum b_i f_i)}{\exp(a_0)} \right) = a_1 - a_0 + \sum b_i f_i$$

اهمیت مدل
این formula

Alternative: $F_i(x) \Rightarrow C(x) \rightarrow \text{horn clause} \rightarrow \text{استاد به سید استراحت}$

Text Classification

page = $\{1, \dots, n\}$

word = $\{\dots\}$

topic = $\{\dots\}$

Topic (page, topic!)

Has Word (page, word)

خود را که با این استناد می کند
اصطلاحی

بر مبنای title دارد

! Topic (p, t)

Has Word (p, w) \Rightarrow Topic (p, t)

حاصل می آید تا ترکیب است

Hypertext Classification

Text Classification

Topic (p, t) \wedge Link (p, p') \Rightarrow Topic (p', t)

Information Retrieval

In Query (word)

Has Word (page, word)

Relevant (page)

In Query (w) \wedge Has Word (p, w) \Rightarrow Relevant (p)

Relevant (p) \wedge Links (p, p') \Rightarrow Relevant (p')

idea from Page Rank

Entity Resolution → Given database, find duplicate records
like google scholar want to merge duplicate citation

Need Transitivity → hard problem in learning
easy for logic

We can resolve the fields as well as records. → two field values refer to same thing

Paper 1: Adam Smith "Pap name" conf1
paper 2: Adam Smith "pap name" conf2 → we can guess that conf1 = conf2

HMM

obs = {obs1, ..., obsN}

state = {st1, ..., stM}

time = {0, T}

State (state!, time) → exactly one state in each time

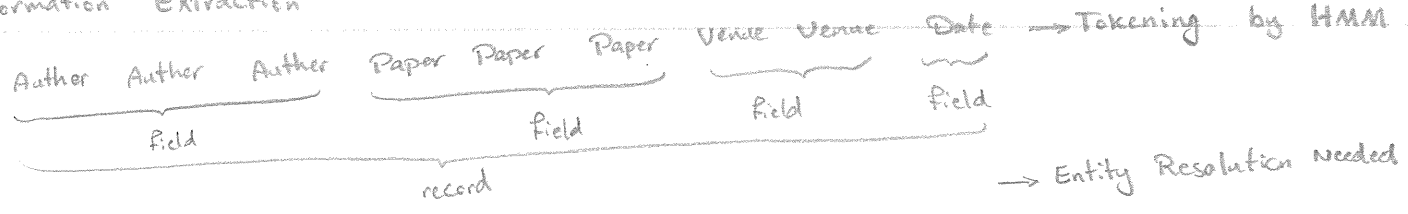
Obs (obs!, time)

State (+S, 0) → each state has different prob. at time 0

State (+s, t) ⇒ State (+s', t+1)

Obs (+o, t) ⇒ State (+s, t)

Information Extraction



Statistical Parsing

Sentence ⇒ Most probable parse

Symbolic Parser + Probability → Statistical Parser

CFG { A → BC
A → a

WCFG
↳ weighted

NP → N Probabilistic CFG { 0.7 NP → N
N → Det N { 0.3 NP → Det N ⇒ Weighted CFG

Query Predicate → John ate the pizza. ⇒ NP (2, 4)
1 2 3 4
↑ NP

Tips

- Not make every predicate closed world → nothing to infer!
 - Conjunctions are closer to probability
 - Complexity → low close arity
low # of constants
short inference chains
- ⇒ go for smallest MLN that work then expand

* Good Design

- Automatic / Semi - Automatic
- Give simple users reasonable results

* Tests available

- Training Accuracy
- Magnitude of Data
- Kernel Analysis → when kernel becomes like Identity Matrix?

↳ what kind of data?

$$\begin{aligned} \Phi(x_1) &= [1, 0, \dots, 0]^T \\ &\dots \\ \Phi(x_c) &= [0, \dots, 0, 1]^T \end{aligned} \Rightarrow \Phi = e^{-\gamma \|x_i - \bar{x}\|} = \begin{cases} i=j, 1 \\ i \neq j, 0 \end{cases}$$

$\Phi = I$

• Overfitting

* Improving Performance

- Data Scaling → overfitting for some kernels can be avoided
- prevent large value features dominance
- Methods → linear
- zero mean, $\Sigma = I$ (whitening)
- Use same scaling methods and factors on training and testing
 - ↳ training and testing are from the same distribution so max, min, var, ... are approx. similar
- depend on kernel use range $[0, 1]$ or $[-1, +1]$ or ...

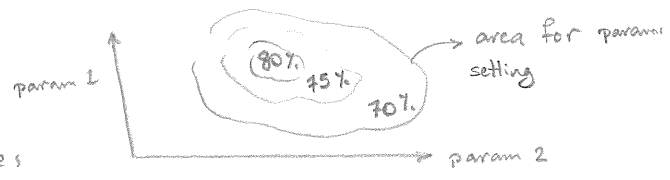
• Parameter Selection → Regularization

- kernel parameters → Polynomial degree
- kernel width
- Cross Validation Parameter → k-fold
- Boost training (?)
- Contour of CV accuracy

only for training set

grid of parameters

↳ try proper ranges



- Selecting kernels → RBF first → good for most data
 - fewer params
 - then Polynomial

• SVM Procedure

1. Data Scaling → linear first
2. Kernel selection → RBF first
3. Use CV to find parameters → γ and C for SVM
4. Use best parameters to train on training set
5. Test

☑ in LIBSVM use easy.py

☑ With bad results when using data scaling → go to test set, print the min and max
 → maybe some min are far from 0 or max are far from 1 → that scaling sucks!!
 → maybe log-scaling

* Machine Learning Software Engineering

- Some people have zero knowledge about ML → not target
- Few people are ML researcher → not target
- Target: general users

↳ automatic / semi-automatic setting

- Multi-class classification → one-against-rest → k binary classifier
- [Hsu and Lin, 2002] → one-vs-one → $\binom{k}{2}$ binary classifier
- ↳ more suitable for kernel SVM
- ↳ user-driven demand! → new direction of research topic!

- User questions → add "one-vs-one" to library
- ↳ add χ^2 kernel → for machine vision simple code

- Extremes → One option for lib → easy to maintain
 ↳ many options for best results → powerful ✓
 ↳ complicated x

Choose Strategy

- People often try all options even if not needed.

↳ some of them has same result but different characteristics.

↳ user have different need ☑ precision / recall

↳ use sample codes to extend lib → ☑ ROC curves
F-score, AUC, ... For CV

- Avoid changes for marginal improvements ⇒ introduce unnecessary complexity

☑ Adding stratified CV → improve stability → worth adding

• shoot for industry → search for lessons learned

- Simplicity vs Performance → sometimes a little less accuracy is meant really simpler algorithm

- Intuition vs. non stable formal methods → visualization and using user intuition
→ simple foolish ideas & Brute force

- Numerical and Optimization Methods → pay attention to algorithm properties

- Scalability
↙
trade time for space
- ☑ keep data in RAM? (Batch)
 - ☑ breakable in smaller problems?
 - ☑ stability for # of training samples
 - ☑ need accurate parameters?
 - ☑ simplification of optimization parts?
 $O(n^2) \rightarrow O(n)$

- Numerical Stability → Catastrophic cancellation → solution, reformulation
→ numerical overflow (floating point) → cases of overflow → reformulation

- legacy → encapsulation
→ documentation
→ learn from questions → users assumption
→ ambiguities

Graph Mining:

(Tsuda - AIST Japan)

* Structural Data

- Nature Tend to use few elements to make complex structures

☑ DNA \rightarrow A, C, G, T

☑ Proteins \rightarrow 20 Amino Acids

- Structure Hidden in Sequences

- infer hidden structures \rightarrow classification

☑ Biological Graph \rightarrow RNA, Gene Expression Data,

☑ Molecular Graph \rightarrow labeled graph \rightarrow nodes: elements

\searrow edges: bonds

- Biological Network

- Networks made of 1000's of nodes / 100,000's of edges

- Possible prediction

Given:	Sequence	Structure
	Structure	\Rightarrow Function
	Expression	Interactions
	Pylogeny: history	Localizations

* Itemset Mining

- Frequent Item set mining \rightarrow all set of elements appearing $\geq \alpha$
 \downarrow minimum support parameter

• Σ set of elements

$0 \leq \alpha \leq T$

• $T \rightarrow$ transactions $\rightarrow t \subseteq \Sigma$

\hookrightarrow usually constant

\searrow $T = \{t_1, \dots, t_m\}$

• appearance: $X \subseteq t$

• occurrence: $occ(X, T) = \{t \in T : X \subseteq t\}$

• frequency: $Fr(X, T) = |occ(X, T)|$

• Frequent: $Fr(X, T) \geq \alpha \rightarrow X$ is Frequent

- Hasse Diagram

- Edge: Adding item

- Need a tree to avoid duplication \rightarrow need a simple tree \rightarrow

BFS, etc.

- Backtracking Algorithm \rightarrow Find Frequent sets \rightarrow FP growth algorithm

- Monotonicity \rightarrow decrease only

- DFS \rightarrow Prune if support $< \alpha$

- Confidence: from data provided to guess the prob. of the rest of the items

☑ $(A, B, C) \rightarrow A, B \Rightarrow C$

- Problem: huge # of frequent itemsets
 - ↳ hard to analyse
 - ↳ many of them are similar
 - ↳ solution: Closed Pattern Mining

* Closed Pattern Mining

• The maximal set among all itemset with same occurrences

• Brute Force Alg

- ↳ generate all freq. sets → test maximality
- exponential size
- delay: time between two pattern outputs
- Brute Force is exponential delay w.r.t pattern size

• To get Linear delay

- Must jump from closed set to closed set
- How? → Reverse Search

• Reverse Search

- Correctness proof mathematically
- Duplication → Naive Backtracking
- Duplication Marking → Exponential Memory
- Visit all nodes without duplication

- Reduction Map

↳ Mapping from children to parent

↳ Given closed set X → shrink till the occurrence change → closure operation
 (A, B, C, D) $(A, B, \cancel{C}, \cancel{D})$ $\bigcap \{t \in T : X \subseteq T\}$

↳ Arrows are in reverse direction

↳ generate all children candidates → apply reduction map → Remove if not coming bac

- Prove the correctness

↳ Reduction maps meet all nodes

↳ Find the root node from any node (by applying reduction map repeatedly)

↳ Children generation is inverse of reduction map

- Algorithm: LCM

↳ prefix preserving closure extension → repeat: add item and taking closure

↳ ensure any closed set is generated from unique parents → choosing the smallest index of matching items → change DAG to Tree

↳ linear delay in pattern size

* Ordered Tree Mining

- Finding all frequent substructures

- Labeled ordered trees

- Rooted
- Ordered
- Labeled

- Frequency of a pattern tree $(T) \rightarrow$ # root occurrences of T in data

- Given: Collection of labeled order tree σ \rightarrow Task: discover Freq. order tree
 Frequency $\geq \sigma$
 without duplicates

- Naive Alg

- Sorting from the smallest tree \rightarrow grow pattern by adding a new node 1 by 1
- exponential ways, need explicit duplication test
- solution: use DFS \rightarrow Rightmost Expansion

\hookrightarrow attaching a new node on the rightmost branch

- Alg.

- Enumerate all freq. order tree by Backtracking
- Tree extension \rightarrow Rightmost Expansion \rightarrow No duplication

* Unordered Tree Mining

- No order between children \rightarrow exponentially many isomorphic trees

- Canonical Representation \rightarrow code $(T) = ((depth_{v_1}, label_{v_1}) \dots (depth_{v_k}, label_{v_k}))$

- special case: ordered tree \rightarrow sort \rightarrow lexicographically maximum code

- Left heavy condition \rightarrow ordertree is canonical \Leftrightarrow code $(T_{v_1}) \geq$ code (T_{v_2})



* Graph Mining

- GSPAN \rightarrow Most popular

- Canonical Representation \rightarrow No root \rightarrow Choose one

- \swarrow DFS order \rightarrow number nodes \rightarrow DFS code

- \searrow Some edges not traversed \rightarrow backward edge (dest# < src#)

- (src, dst, src_label, edge_label, dest_label)

- Multiple DFS code \rightarrow different { start point, child order } \rightarrow Minimum DFS code (sort)

- Reduction Map

- Removing the tail of min DFS code \rightarrow preserve minimality

• Children Generation

- candidate: add element to DFS code → candidate minimum? $\xrightarrow{\text{No}}$ remove it

• Minimality Check

- Reconstruct graph from DFS code

- Derive min DFS code by trying all DFSs → speed up by traversing min label only

- minimal code \neq original? → prune it.

• gSpan → OK

cloSpan → not ok → [Takigawa et. al, 2011]

* Dense Module Enumeration

• Challenges

- False negative edges → go beyond clique search

- False positive edge → assign confidence scores to edges

• Previous

- Clique percolation

- Partitioning

- Heuristic local search

• Exhaustively enumerate all dense subgraph

- Efficient

- Detect overlaps

- Min density for outgoing module

- All modules output → satisfying threshold

• Problem

- Interaction net $G = (V, E)$

- Edge weight $w \in [0, 1]$

- Density
$$d(U) = \frac{\sum_{\{u,v\} \in E} w(\{u,v\})}{|U|(|U|-1)}$$

⇒ Find $U \subset V$
$$\begin{cases} d(U) \geq \theta \\ c_1(U) \geq nt \rightarrow \text{criteria 1} \\ c_0(U) \geq no \rightarrow \text{criteria 0} \end{cases}$$

• Enumeration

- set up a search tree → pruning by anti-monotonicity
↳ $\text{freq}(\text{entity}) < \text{freq}(\text{sub entity})$

- density is not monotonic → pruning impossible
for lexicographical ordering

- need a search tree with monotonically decreasing density → Use reverse search

• Reverse Search

- Specify search tree → use reduction map
↳ generate a parent from a child → remove the node with smallest degree → always density ↑

• Alg For Enumeration

- 1. set of children from a parent
- 2. choose children satisfying reduction mapping
- 3. Prune if no children remain

* Graph ML Examples

• Given set of graph \implies cluster them in groups

• Graph Regression

• Substructure Representation

- 0/1 vector of pattern indicator

↳ huge dimensionality

↳ need graph mining for selecting features

• Frequent Substructure Mining

- Enumerate all patterns in at least m graph

- Gspan

• Discriminative patterns

- $w_i > 0 \rightarrow +$ class

$w_i < 0 \rightarrow -$ class

- weighted Substructure Mining

- Patterns with large freq difference

• Multi class

- $w_{e_i} > 0 \rightarrow$ graph $i \in$ class l

- search patterns overrepresented in a class

* EM-Based graph clustering

• learning feature & model at the same time

• L1 regularization \rightarrow make graph sparse

• E-step \rightarrow Mining step \rightarrow M-step

• Prob.

- Binomial Mixture

$$P(x|\theta) = \sum_{e=1}^c \alpha_e P_e(x|\theta_e)$$

↖ weight
↖ parameter for cluster e

$$P_e(x|\theta_e) = \prod_k \frac{\exp(\theta_{ek} \cdot x_k)}{1 + \exp(\theta_{ek})}$$

• EM

$$\operatorname{argmax}_{\theta} \sum_{i=1}^n \log \sum_{e=1}^c \alpha_e P_e(x_i|\theta_e)$$

E-step: $r_{ei} = P(y=e | x_i)$

M-step: estimate θ_e

\implies computationally $\bullet \rightarrow \dots$

• Use L1-regularized log likelihood

• E-step

- Active Pattern $F = \{k \mid \exists e : \theta_{ek} \neq \theta_{ok}\}$

- E-step computation for active patterns ✓

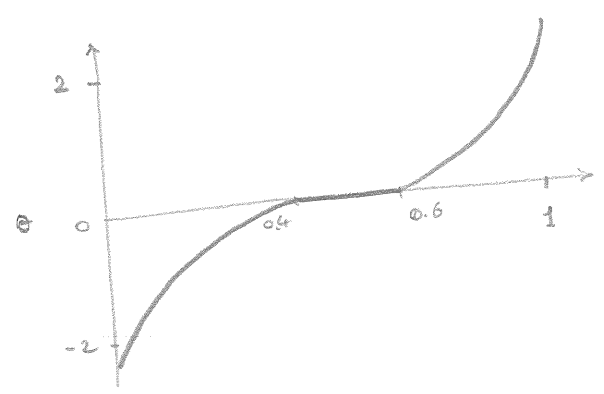
$$P(y=e|x) = \frac{\alpha_e \prod_{k \in F} P_{ek}(x_k | \theta_{ek})}{\sum_e \dots}$$

• M-step

- use graph mining to find active patterns

$$\theta_{ek} = \begin{cases} \log \frac{n_{ek} - \lambda_e}{1 - (n_{ek} - \lambda_e)} \\ \theta_{ok} \\ \log \frac{n_{ek} + \lambda_e}{1 - (n_{ek} - \lambda_e)} \end{cases}$$

λ_e = weighted occurrence
 n_{ek} = occurrence prob in cluster
 n_{ok} = occurrence prob overall



For active pattern, occurrence prob is far from avg.

* Graph Boosting

• Given : graph \rightarrow predict class label (+/-)

Typically features are given

• Motivation \rightarrow Missing Descriptors

\rightarrow New features \rightarrow informative patterns \rightarrow greedy pattern discovery (Boosting + gSpan)

• Linear Programming Boosting

\rightarrow less calls to graph mining
 \rightarrow faster than AdaBoost

• Formulation

$$\min_{\alpha, \epsilon} \|\alpha\|_1 + C \sum_{n=1}^e \epsilon_n$$

$$\text{s.t. } y_n \alpha^T x_n \geq 1 - \epsilon_n, \epsilon_n \geq 0$$

- sum of hinge loss and L1 regularizer

- dual problem : \uparrow # of weight variables $\Rightarrow \uparrow$ # of constraints

- Finding the most violated constraint \rightarrow search by weighted substructure mining

* Entire Regularization Path

• LARS - LASSO, LISVM

• Start from small set of features \rightarrow trace solution trajectory of L1 reg. learning \rightarrow add more features

• LASSO Regression

$$B(\lambda) = \underset{\beta}{\operatorname{argmin}} L(y, X\beta) + \lambda \|\beta\|$$

$$\lambda \in \{0.001, 0.01, 0.1, 1, 10, \dots\}$$

A: active feature set

- Piecewise Linear Path \rightarrow At turning point \rightarrow add new feature to A
- Cross validation \rightarrow Grid Search \rightarrow Solve QP many times x .
 - \rightarrow delete feature from A
 - \rightarrow Path following \rightarrow does not include QP $\rightarrow \lambda_{cv}^*$

• Feature Space of Patterns

- Representation $\alpha_i = (\alpha_{it})_{t \in T}$, $\alpha_{it} \in I$ ($t \subseteq G_i$)
- Find pattern t that minimize $d_t = \text{step size}$
- Tree shaped search space \rightarrow node: pattern
 - \rightarrow edge: add edge in each level

* kernels

- learning \rightarrow general purpose learning \rightarrow SVM, PCA
 - \rightarrow problem specific kernel
- kernel \rightarrow mapping from a space to a linear (or simple) space
 - \rightarrow show similarity between two things

- implicit mapping $\Phi \rightarrow$ similarity measure $k \rightarrow k(\alpha, \gamma) = \Phi(\alpha) \cdot \Phi(\gamma)$

- similarity measure \rightarrow invariance or other a priori knowledge
 - \rightarrow in vision \rightarrow invariance to rotation
 - \rightarrow class of function
 - \rightarrow possibility infinite dimension \rightarrow still computationally efficient

- Kernel Trick $\rightarrow \alpha \cdot \gamma \Rightarrow \Phi(\alpha) \cdot \Phi(\gamma)$
 \rightarrow graph kernel

\rightarrow distance based: $d(\alpha, \gamma) \Rightarrow k(\alpha, \alpha) + k(\gamma, \gamma) - 2k(\alpha, \gamma)$

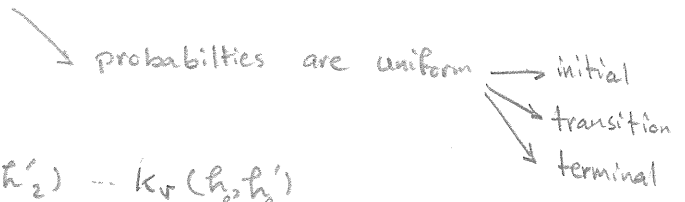
- Valid kernels \rightarrow Positive definite $\int K(\alpha, \gamma) \cdot d\alpha \cdot d\gamma \geq 0$
 \rightarrow eigen values are non-negative

- Design kernels \rightarrow Convolution
- \rightarrow Generative Model-based

* Marginalized kernel

- define kernel on the hidden variable
- Joint kernel: count kernel with context info \rightarrow marginalize Joint k .
- Marginalized graph kernels
 - both vertices and edges are labeled

- Label path $\{v_1, e_1, v_2, e_2, \dots, v_n\}$ \rightarrow by random walk



- kernel $k(h, h') = \begin{cases} 0 \\ k_v(h_1, h'_1) k_e(h_2, h'_2) \dots k_r(h_n, h'_n) \end{cases}$

\hookrightarrow Marginalized $k(G, G') = \sum_h \sum_{h'} p(h|G) p(h'|G') k(h, h')$

✓ ECoG

* Invasive vs. noninvasive BCI

- BCI: human intention \Rightarrow technical control signal
 - \hookrightarrow without using activity of muscles / peripheral nerves
 - \hookrightarrow Brain Signal \Rightarrow Decoding \Rightarrow Control Device

✓ BBCI Pong \rightarrow Subject: assume two string to control the cursor
 after 30 min \rightarrow no need to string pulling
 skill acquisition

• Challenges \rightarrow cerebral cocktail party problem \rightarrow different mental activity

Brain Rhythms \rightarrow α rhythm in visual cortex \rightarrow overlapping things
 \rightarrow μ rhythm in sensory cortex \rightarrow having same rhythm

Variations \rightarrow single-trial vs. averaging

\rightarrow season to season variability \rightarrow day 2 day diff.

\rightarrow inter subject variability

non stationarity shifting distribution

• Data \rightarrow channel + time \rightarrow for imagining of brain state

* Common Spatial Pattern Analysis

• 2 class: Data \rightarrow whitening of joint matrix \rightarrow Rotation (R matrix should be chosen so one matrix should be diagonal and the other should be 1-Diagonal).

* In Mhw non-stationarity is neglected.

* Bad Cross Validation

• Test set should not be used.

• Preprocessing of data regarding all data

• Features are selected on whole data \rightarrow ✓ automatic feature selector

• Parameter selection from all data

• Insufficient validation for paradigms with block design

- Block design \rightarrow blocks are dependant

\rightarrow leave-one-block-out validation \rightarrow no more optimistic statements

* Non Stationarity

• Mathematical Flavors

- Bias adaptation
 - Covariance shift \rightarrow we have model of world, then
 - Projection method \rightarrow project data to lower data that is stationary in that
 - Mixed effect model \rightarrow subject dependence
 - Co-adaption \rightarrow Machine and subject adapt to each other
- \searrow Non stationarity in R^h

• Weighted linear Regression for Covariance shift compensation

- $P_{train}(x) \neq P_{test}(x)$
- use weight in MSE \rightarrow unbiased it again \rightarrow have proof
- density estimation in high dimension is hard \rightarrow weights make it nightmare

• Projection Technique: Splitting into stationary and substationary subspace (SSA)

- PCA: $X =$ orthogonal mixing \times uncorrelated sources

$$X = A \begin{bmatrix} S_1 \\ \vdots \\ S_d \end{bmatrix} \begin{matrix} \rightarrow \text{max var} \\ \rightarrow \text{min var} \end{matrix}$$

- ICA: $X =$ arbitrary mixing \times independent sources

$$X = A \begin{bmatrix} S_1 \\ \vdots \\ S_d \end{bmatrix}$$

JMLR website

- SSA = Stationary Subspace Analysis

$$X_E = A \begin{bmatrix} S_d^s \\ S_d^h \end{bmatrix} \begin{matrix} \rightarrow \text{stationary sources} \\ \rightarrow \text{non stationary sources} \end{matrix}$$

to get a projection that everything stationary

$$X(t) = A \cdot S(t) = \begin{bmatrix} A^s & A^h \end{bmatrix} \begin{bmatrix} S^s(t) \\ S^h(t) \end{bmatrix}$$

• SSA

objective

$$\begin{bmatrix} \hat{S}^s(t) \\ \hat{S}^h(t) \end{bmatrix} = \begin{bmatrix} \hat{B}^s \\ \hat{B}^h \end{bmatrix} x(t)$$

$$\hat{B}^s = \underset{RR^T = I}{\operatorname{argmin}} \sum_{i=1}^N (-\log \det \sum_{i=1}^N S_i + \lambda_i^s \lambda_i^h)$$

Computationally complex

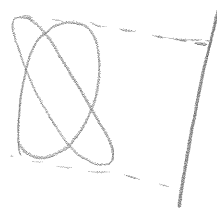
- Optimizing

$$B^{new} \leftarrow R B^{old} \Rightarrow \text{Find } R \text{ to ease calculation of gradients}$$

under these rotation they look the same

we need a third one

how many rotation?



$$N > \frac{D-d}{2} + 2. \rightarrow \text{non-stationarity needs both mean and covariance}$$

\Rightarrow Paper: Kraljic et al. 2012 JMLR

\Rightarrow BBCI.de

\Rightarrow paper: Somak et al. 2012 JNE

* Bias in computing Covariance

- sample size large

$$\mu = \frac{1}{n} \sum \alpha_k$$

$$\Sigma = \frac{1}{n-1} \sum (\alpha_k - \mu)(\alpha_k - \mu)^T$$

- sample size small \rightarrow biased!

□ BCI, Financial Analysis

- large eig values \rightarrow too large

- small eig values \rightarrow too small

- using shrinkage estimation

• Shrinkage

$$\tilde{\Sigma}(\gamma) = (1-\gamma) \hat{\Sigma} + \gamma \nu I$$

ν = average eigen values \rightarrow eigen values change toward it.

$$\tilde{\Sigma} = (1-\gamma) \hat{\Sigma} + \gamma \nu I = (1-\gamma) V D V^T + \gamma \nu I = V ((1-\gamma) D + \gamma \nu I) V^T$$

$\hookrightarrow \gamma = 0 \rightarrow$ no shrinkage

$\hookrightarrow \gamma = 1 \rightarrow$ spherical covariance matrix

- optimal selection of γ

$$\gamma^* = \frac{n}{(n-1)^2} \frac{\sum_{i,j,k} \text{var}_k (Z_{ij}(k))}{\sum_i (s_{ii} - \nu)^2 + \sum_{i \neq j} s_{ij}^2}$$

$$\hat{\mu} = \begin{bmatrix} \mu_1 \\ \vdots \\ \mu_d \end{bmatrix}$$

$$\hat{\Sigma} = \begin{bmatrix} s_{11} & & \\ & \ddots & \\ & & s_{jj} \end{bmatrix}$$

$$Z_{ij} = \begin{pmatrix} (\alpha_k)_i - (\hat{\mu})_i \\ (\alpha_k)_j - (\hat{\mu})_j \end{pmatrix}$$

* ELG

- Presurgical localization of brain part causing epilepsy \Rightarrow [Schalk]
- good to understand human

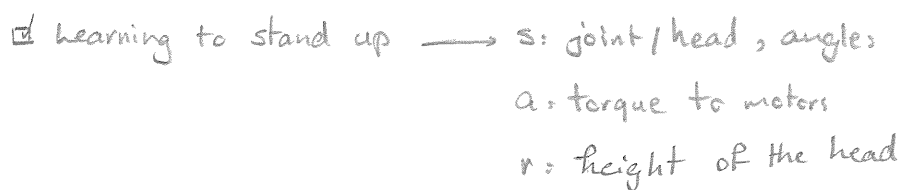
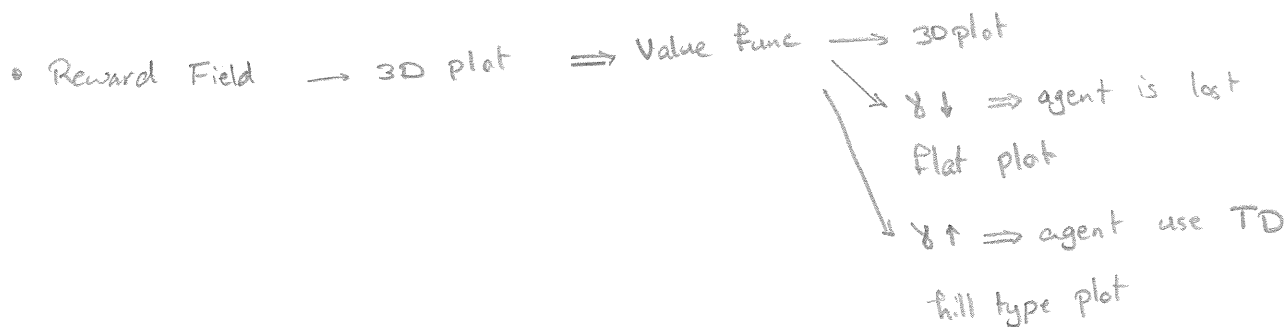
□ good for experimenting \rightarrow Index - vs - rest, Thumbs - vs - rest for finger movement

* Motory Decision

- 600 ms to action \rightarrow start of process
- 300 ms to action \rightarrow peak signal
- 150 ms to action \rightarrow the point of no return from the action

* RL

- Value Func: expected future reward
- learn action policy: $s \rightarrow a$ to maximize reward
- Temporal Difference: diff between what you expected & what you have

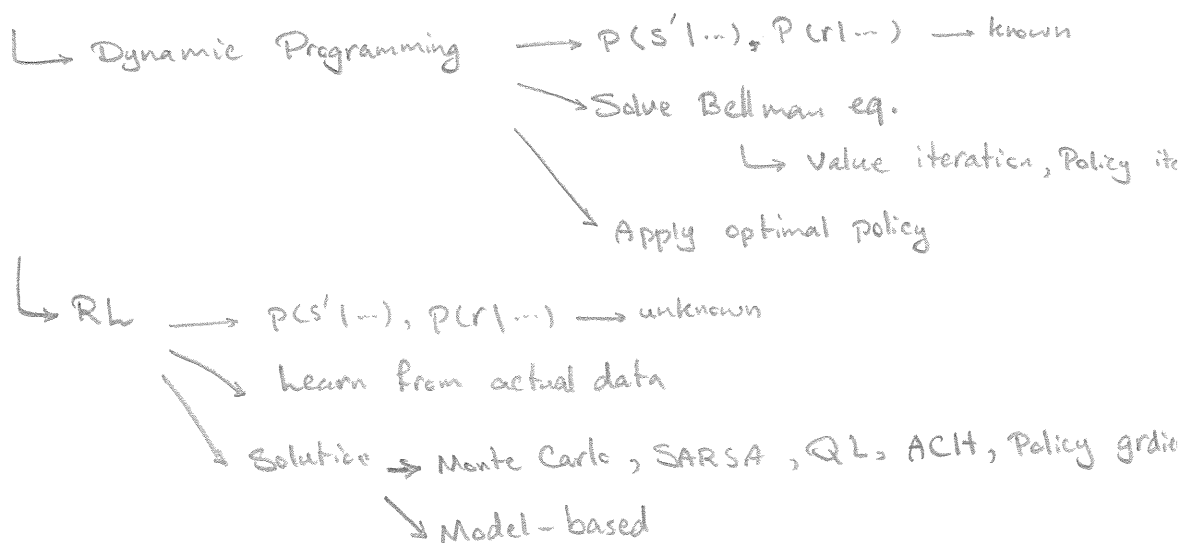


• MDP

- optimize policy: max cumulative reward

- \rightarrow finite horizon
- \rightarrow infinite horizon ($\gamma \leq 1$)
- \rightarrow average reward

- Solutions



* Actor-Critic Learning

- Actor : parameterize policy
- Critic : learn value function

- in table or neural network (function approximator)

• TD error $\rightarrow \delta$: diff of actual value and predicted \rightarrow supervised

• Update \rightarrow critic : $\Delta V = \alpha \delta$

\rightarrow actor : $\Delta \omega = \delta \frac{\partial P(a|...)}{\partial \omega} \rightarrow$ reinforce action by δ

* SARSA and QL

• Action-Value function: $Q(s, a) = \mathbb{E}[r(t) + \gamma r(t+1) + \dots]$

• Action selected

- ϵ -greedy $\rightarrow \arg \max_a Q(s, a)$ with prob $1 - \epsilon$

\rightarrow all the suboptimal actions takes similar weights

- Boltzman \rightarrow action with higher $Q(s, a)$ is taken with higher prob.

$\rightarrow \beta$: parameter controlling greediness

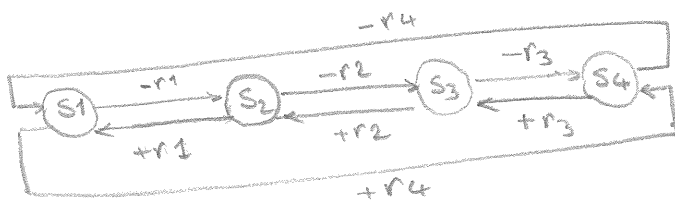
• SARSA : on-policy update

\rightarrow take your decision based on policy \rightarrow then update $Q(s, a)$

• QL : off-policy update

\rightarrow choose next action based on greedy decision \rightarrow the update $Q(s, a)$

□ "Loss to Gain" task



$\frac{r_2}{r_1} \rightarrow$ ratio important

γ small \rightarrow oscillation

β large \rightarrow oscillation \rightarrow low explore

* Issues of RL



• Define state, action

• Design reward \rightarrow balance between reward & punishment

• Which algorithm

• How to set meta parameters \rightarrow learning rate, exploration

• Partial Observability \rightarrow QL have problem, need hidden states

• Non stationary \rightarrow change of environment before convergence

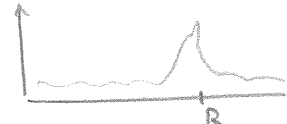
• Modules & Hierarchy

• Dopamine neurons code TD error [Sculitz et al. 1997]

☑ dopamine neurons

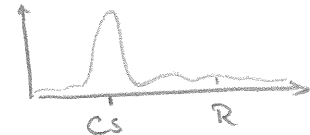


→ unpredicted reward



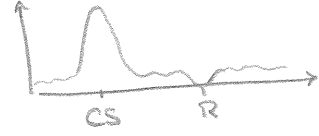
⇒

→ predicted reward



⇒

→ omitted reward



- اگر نماند جایزه نه برنت ترشح می کند
- اگر منتظر باشد زمان انتظار ترشح می کند
- اگر منتظر باشد ولی پاداش ندهیم، دیرزمان پاداش هم ترشح می کند

• Basal Ganglia For RL? [Doya 2007] [Doya 2000]

- state & action → cerebral cortex
- reward pred → striatum $Q(s, a), V(s)$
- temporal diff → dopamine neurons δ
- action selection → pallidum
- memory (?) → thalamus

☑ addicts → problem with dopamine release

Hebbian Learning → dopamine dependant plasticity

☑ Monkey Free choice task [Samejima et al. 2005]

- ↳ dissociated action & reward
- ↳ neurons in striatum is recorded → action-value coding in striatum

☑ Forced and Free Choice Task → Iko, OIST

↳ Generalized Q-learning Model [Ito & Doya, 2009]

↳ α_1 : learning rate α_2 : forgetting rate k_1 : reinforcement k_2 : aversio

- Model Fitting by particle filtering

• Neural activity in the Striatum

• Hierachy in Cortico-striatal Net [Voorn et al 2004]

- Dorsolateral → motor
 - ↳ early action → what action
- Dorso medial → frontal
 - ↳ action value → what context?
- Ventral → limbic
 - ↳ state value → whatever worth doing?

* Temporal Discount Factor

- γ large \rightarrow go for far goals \rightarrow impulsive
- γ small \rightarrow go for only near goals \rightarrow depression
- control the character of agent
- γ should be large enough \rightarrow too large : slow learning
- Neuro modulators for Metalearning [Doya, 2002]
 - to tune parameter
 - Dopamine : TD error
 - Serotonine : higher firing during waiting for reward , stops before giving up
 - \hookrightarrow regulation for patient
 -
 -
- \square Delay Tone - Food - Tone - water Task

* Intro

- minimize a function subject to some constraints
- global min is hard to find
- special case: Linear Programming \rightarrow Famous \rightarrow polynomial time
- many problems can be solved \rightarrow has a condition: convex
 - \rightarrow provide heuristics for non-convex problem
- new application since 1990's
- advances \rightarrow interior-point method
 - \rightarrow first-order algorithm

* Convex Set

• $x_1, x_2 \in C, 0 \leq \theta \leq 1 \implies \theta x_1 + (1-\theta)x_2 \in C$

بر خط بین دو نقطه شکل اردن
شکل است

• Basic sets

- Affine set
- Half space
- Polyhedron \rightarrow finite
- Ellipsoid \rightarrow quadratic, positive definite
- Norm ball \rightarrow
- Positive Semidefinite Cone $\rightarrow S_+^n$

• Intersection of convex sets is convex \rightarrow using definition of convexity

* Convex Function

• $(x, f(x))$ \implies بر خط بین دو نقطه، تابع دایره \implies Jensen Inequality
 $(y, f(y))$ \implies طرف تابع نه می برد. $f(\theta x + (1-\theta)y) \leq \theta f(x) + (1-\theta)f(y)$

• Basic Examples

- linear combination
- $e^x, -\log x, x \log x$
- $x^\alpha, \alpha \geq 1, \alpha \in \mathbb{R}^+$
- f convex $\implies -f$ concave

• Epigraph $\rightarrow \text{epi}(f) : \text{تقاطع نصفه}$
 $\searrow f : \text{convex func} \Leftrightarrow \text{epi}(f) \Leftrightarrow \text{convex set}$

• Differentiable Conv. Func
 \searrow well defined min \rightarrow convexity is good for minimization

- Showing Convexity
 - \hookrightarrow 1. Def
 - \hookrightarrow 2. $\nabla^2 f(x) \succcurlyeq 0$
 - \hookrightarrow 3. a function of convex function that preserve convexity

□ point wise max, point wise supremum (infinite # for taking max)

• Conjugate Function

- $f^*(y) = \sup_{x \in \text{dom } f} (y^T x - f(x)) \rightarrow f(x), y^T x \text{ تقاطع نصفه پیشنهاد}$

- $f^* \rightarrow \text{convex} \rightarrow$ not matter what is f

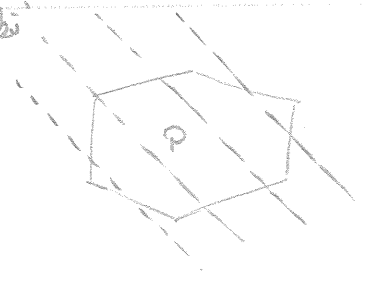
$\searrow \text{sup (convex)} = \text{convex} \rightarrow \text{fix } x : \text{linear} : \text{convex}$

- good for duality

* Linear Programming

تقاطع نصفه P تقاطع خطی

subject to $\left\{ \begin{array}{l} \min \quad c^T x + d \\ \text{s.t.} \quad Gx \leq h \\ \quad \quad Ax = b \end{array} \right.$



• look at the lines, move them, to find the result

• usually found in piece-wise linear minimization

- $\min f(x) = \max (a_i^T x + b_i)$

- transform to $\min t$
 s.t. $a_i^T x + b_i \leq t, \forall i$



• ℓ_1 -norm approximation

- def $\|y\|_1 = \sum_k |y_k|$

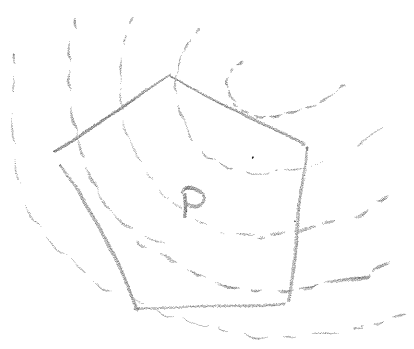
- min $\|Ax - b\|_1$

- min $\sum y_i$
 s.t. $-y \leq Ax - b \leq y$

• ℓ_∞ -norm ($\|y\|_\infty = \max_k |y_k|$)

- min $\|Ax - b\|_\infty$

- min y
 s.t. $-y\mathbf{1} \leq Ax - b \leq y\mathbf{1}$



* Quadratic Programming

• min $(\frac{1}{2}) x^T P x + q^T x + r$
 s.t. $Gx \leq h$

• Linear program with random cost

min $c^T x$
 s.t. $Gx \leq h$

$c \sim N(\bar{c}, \Sigma)$

• Kernel Formulation

- min $f(xa) + \|a\|_2^2$ \rightarrow Euclidean dist

- dual \rightarrow min $f(y) + y^T Q^{-1} y$
 $Q = x x^T \rightarrow$ kernel matrix

* Geometric Programming

• $f(x) = \sum_{k=1}^K c_k x_1^{a_{1k}} x_2^{a_{2k}}$

min $f_0(x)$
 s.t. $f_i(x) \leq 1 \quad \forall i$

• Second-order Cone Programming (SOCP)

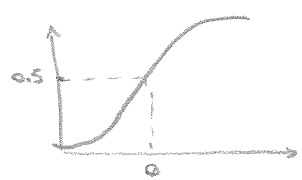
min $f^T x$
 s.t. $\|A_i x + b_i\|_2 \leq c_i^T x + d_i$

• Robust linear program (stochastic)

- min $c^T x$
 s.t. $\text{prob}(a_i^T x \leq b_i) \geq \eta \rightarrow \eta = 90\%$

$\eta \uparrow \Rightarrow$ stay away from edges of polyhedron
 $\eta \downarrow \Rightarrow$ more toward non-convexity \Rightarrow sharp edge $\approx \eta = 50\%$

- convert to LP (require $\eta \geq 50\%$)



$\Phi^{-1}(\eta) > 0 \Leftrightarrow \eta \geq 0.5$
 convex

Φ : cumulative density function

* Semidefinite Program (SDP)

- $\min c^T x$
s.t. $x_1 A_1 + x_2 A_2 + \dots + x_n A_n \preceq B$
- include many non-linear constraints as special cases
- semidefinite relaxation
 - to solve non-convex \rightarrow find bounds
 - \searrow heuristic
 - reformulation trick \rightarrow change the variable \rightarrow try to make constraints convex
 - \searrow expand formula
 - replace constraints with weaker ones \rightarrow relax!
 - \searrow for example to obtain lower bound

* Conic Optimization

- Inequality \rightarrow constraint $x \in K \rightarrow K$: convex cone
- Proper Cones \rightarrow closed, pointed, non-empty interior
- Cone linear Program
 - $\min c^T x$
s.t. $Ax \preceq_K b$
 - popular \rightarrow modeling: small # of primitive cones is sufficient
 - \searrow algorithm: extend LP

• Second order cone program

- $\min c^T x$
s.t. $\|B_{k0} x + d_{k0}\|_2 \leq B_{k1} x + d_{k1}$

- $R_+^n, Q^p, S^p \rightarrow$ enough cone set for many applications
 - \hookrightarrow exponential & logarithm

* Duality

• Primal	$\min c^T x$ s.t. $Ax \preceq b$	\iff	Dual	$\min -b^T z$ s.t. $A^T z + c = 0$ $z \succeq 0$
	\hookrightarrow optimal p^*			\hookrightarrow optimal d^*

- $p^* \geq d^* \rightarrow$ weak duality
- $p^* = d^* \rightarrow$ strong duality

• $x \succeq_* b \iff x \succeq_{K^*} b \rightarrow$ using proper cone K

* Boosting & (Schapire - Princeton U.)

- To finding single highly accurate prediction rule
 - ↳ From "rule of thumbs" that are often correct
 - ↳ how to choose examples in each round of finding rule of thumbs
 - ↳ hardest examples
 - ↳ how to combine those rules
 - ↳ weighted majority vote of them
- have weak learning algorithms with accuracy better than random \rightarrow acc $>$ 55% (rule of thumbs)
 - \Rightarrow make very high accuracy classifier (e.g. 99%)

* History

- Arise from PAC model : From weak PAC to strong PAC
 - ↓
 - IF procedure found we can reach perfect accuracy
 - ↓
 - learning is in two states \rightarrow completely learn
 - \rightarrow random performance
- First boosting algorithms
- Ada Boost \rightarrow ...

* Ada Boost

1. m training example $(x_1, y_1) \dots (x_m, y_m) \rightarrow y_i \in \{0, 1\}$
2. T rounds
 - ↳ construct dist D_t on $\{1, \dots, m\} \rightarrow$ weights $\rightarrow \begin{matrix} \in \mathbb{R}^+ \\ \sum = 1 \end{matrix}$
 - ↳ Find weak classifier h_t
 - ↳ measure error ϵ_t on $D_t \rightarrow$ Prob that an example misclassified
3. Combine classifiers H_{final}

- Construct D_t
 - $D_1(i) = 1/m$
 - $D_{t+1}(i) = \frac{D_t(i)}{Z_t} \times \begin{cases} e^{-\alpha_t} & y_i = h_t(x_i) \\ e^{\alpha_t} & y_i \neq h_t(x_i) \end{cases} \rightarrow Z: \text{normalization factor}$
 - $= \frac{D_t(i)}{Z_t} \exp(-\alpha_t y_i h_t(x_i)) \rightarrow$ increase the weight on misclassified
 - ↓
 - finding hardest examples

• Construct Final Classifier

$$H_{\text{final}}(x) = \text{Sign} \left(\sum \alpha_t h_t(x) \right)$$

• Calculate α

$$\alpha_t = \frac{1}{2} \ln \left(\frac{1 - \epsilon_t}{\epsilon_t} \right) > 0$$

• The training error is always decreasing during the training

$$\hookrightarrow \text{edge } \gamma_t \rightarrow \epsilon_t = \frac{1}{2} - \gamma_t$$

$$\text{training error } (H_{\text{final}}) \leq \prod_t [2\sqrt{\epsilon_t(1-\epsilon_t)}] \leq \prod_t \sqrt{1-4\gamma_t^2} \leq \exp(-2 \sum_t \gamma_t^2)$$

$$\hookrightarrow \text{even if } \gamma \text{ is small } \text{error}(H_{\text{final}}) \leq \exp(-2\gamma^2 T) \quad \gamma_t \gg \gamma > 0$$

\hookrightarrow does not need to know γ or T a priori \rightarrow adaptive boosting

• Proof

$$1. D_{t+1}(i) = D_t(i) + \dots \Rightarrow D_{\text{final}}(i) = \frac{1}{m} \frac{\exp(-y_i \sum \alpha_t h_t(x_i))}{\prod_i Z_t}$$

$$2. \text{training error} \leq \prod_t Z_t$$

$$\text{training error} = \frac{1}{m} \sum \mathbb{I}(y_i \neq H_{\text{final}}(x_i)) = \frac{1}{m} \sum \mathbb{I}(y_i F(x_i) \leq 0)$$

$$\leq \frac{1}{m} \sum \exp(-y_i F(x_i)) = \underbrace{\sum_i D_{\text{final}}(i)}_1 \prod_t Z_t = \prod_t Z_t$$

$$3. Z_t = 2\sqrt{\epsilon_t(1-\epsilon_t)}$$

$$Z_t = \sum D_t(i) \exp(-\alpha_t y_i h_t(x_i)) = \sum_{i: y_i \neq h_t} D_t(i) e^{\alpha_t} + \sum_{i: y_i = h_t} D_t(i) e^{-\alpha_t}$$

$$= \epsilon_t e^{\alpha_t} + (1-\epsilon_t) e^{-\alpha_t} = 2\sqrt{\epsilon_t(1-\epsilon_t)} \quad \hookrightarrow \text{where } \alpha_t \text{ come from}$$

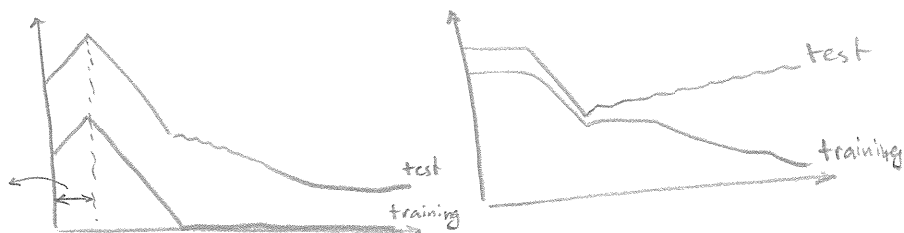
• The test error

- increase in beginning \rightarrow do better

- decrease after round $k \rightarrow$ adding boosting a new classifier each round

\rightarrow overfitting \rightarrow according to Occam's razor

1 classifier \rightarrow dictatorship
2 classifiers \Rightarrow one always wins



training err = 0
test error decreasing \Rightarrow Occam's razor!

• Margin explanation

- idea → training err = validating classification

→ we need away of measuring confidence

- for voting → margins show the winner and confidence

→ AdaBoost is a weighted majority vote

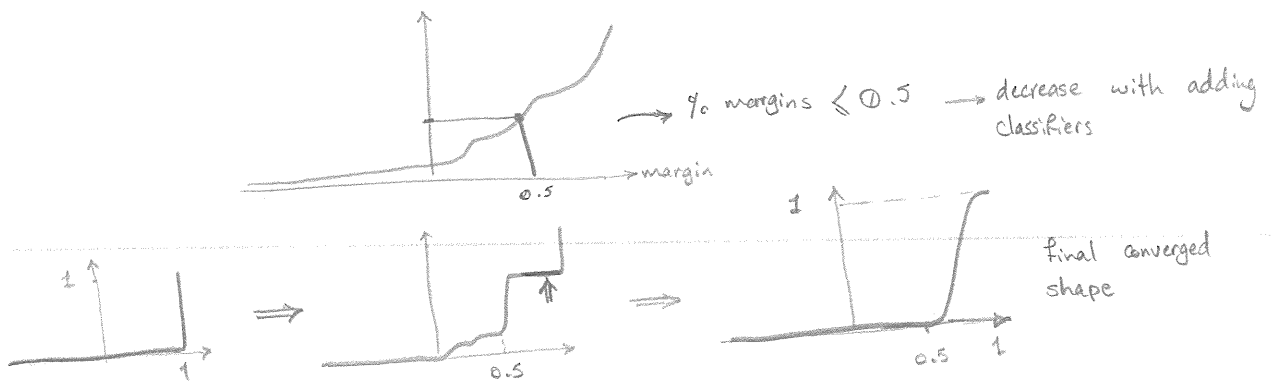
margin = strength of the vote

= weighted fraction - weighted fraction
voted correctly voted incorrectly

↳ sign (margin) → voting correctly (in correctly)

↳ |margin| → confidence

- margin dist = cumulative dist of margins of training example



- larger margins ⇒ better bounds on generalization err

↳ independent of # of rounds

- boosting tends to increase margins of training examples (given weak classifier assumption)

- final classifier → getting bigger

→ getting more look like simpler classifiers

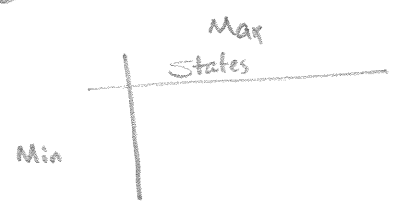
→ larger margins

→ resistant to overfitting → small edge: under fit

→ overly complex weak classifier: overfit

* Game Theory Perspective

- Boosting as a game
- Game



→ Min: dist P over rows
 Max: dist Q over columns

$$\text{Min Loss} = \sum P_{i,j} M(i,j) Q(j) = P^T M Q = M(P,Q)$$

- Sequential Game: Min play before Max

$$\hookrightarrow L(P) = \max_Q M(P,Q)$$

$$\hookrightarrow \min_P L(P) = \min_P \max_Q M(P,Q)$$

→ playing better is usually looks better but in real it makes no difference → Von Neumann's theorem of Minimax

$$\min_P \max_Q M(P,Q) = \max_Q \min_P M(P,Q) = V$$

☑ in Rock-Paper-Scissors ⇒ play with uniform dist on action

- Classic theory weakness

↳ game unknown

↳ game might be extremely large

↳ not adversarial opponent → room for better performance

• Repeated Play

- possible to learn to play well → learn game matrix / opponent

- M unknown

in a loop

- Min chooses P_t
- Max chooses Q_t
- Loss of Min = $M(P_t, Q_t)$
- Min observe loss $M(i, Q_t)$ of each pure strategy i

- Multiplicative-weight Algorithm

↳ punish weight of strategies suffering the most loss

$$P_{t+1}(i) = \frac{P_{t,i} \exp(-\eta M(i, Q_t))}{\text{normalization}}$$

↳ goal: actual avg. loss \leq best average loss + Δ_T

$$\Delta_T = O\left(\sqrt{\frac{\ln m}{T}}\right) \rightarrow \text{regret} \quad \left\{ \begin{array}{l} \text{ind} (\neq \text{columns}) \rightarrow \text{opponent} \\ \text{dep} (\neq \text{rows}) \end{array} \right.$$

- Solving a game

Min play using MW

Max choose best response

• Boosting as a game

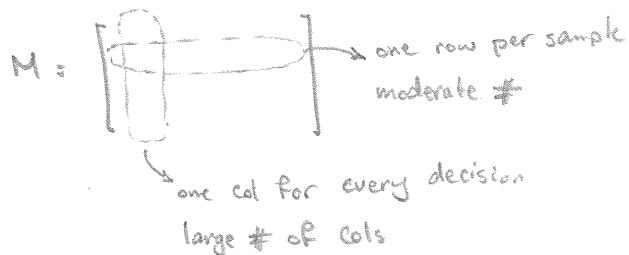
- Min \rightarrow booster (row)

- Max \rightarrow weak learner (col)

- Matrix M

\hookrightarrow row : training exemplar

\hookrightarrow col : weak classifier



$\hookrightarrow M(i,j)$ = which classifier are correct/incorrect on each example

- γ -weak learning assumption

for each dist on example there's weak classifier with weighted err $\leq \frac{1}{2} - \gamma$

\hookrightarrow value (M) $\geq \frac{1}{2} - \gamma$

\hookrightarrow there's a weighted majority classifier that correctly classifies all training examples by margin $\geq 2\gamma$

\hookrightarrow Use MW to solve the game \Rightarrow variant of AdaBoost

* Loss Minimization

• Fitness measuring \rightarrow least square, ...

• why? \rightarrow convergence proof

\searrow decoupling alg from objective

• AdaBoost

- training err (H_T) $\leq \prod_t Z_t$

- $Z_t = E_t e^{\alpha_t} + (1 - E_t) e^{-\alpha_t} = 2 \sqrt{E_t(1 - E_t)}$

$\hookrightarrow \alpha_t$ minimize Z_t

$\hookrightarrow h_t$ minimize $E_t \Rightarrow$ minimize Z_t

- greedy procedure for minimizing exponential loss

$\hookrightarrow \prod_t Z_t = \frac{1}{m} \sum_i \exp(-y_i F(x_i))$

$\hookrightarrow F(x) = \sum_t \alpha_t h_t(x)$

- why exp. loss

↳ encourage exponentially to $F(x_i)$ and y_i have same signs

↳ upper bound on training err

↳ non-convex

- Coordinate Descent

$$F(x) = \sum_t \alpha_t h_t(x) = \sum_j \lambda_j g_j(x)$$

↳ find λ_j to minimize loss

$$L(\lambda_1, \dots, \lambda_N) = \sum_i \exp\left(-y_i \sum_j \lambda_j g_j(x_i)\right)$$

- Ada Boost as coordinate descent → powerful over huge space of functions

↳ choose coord. $\lambda_i \Rightarrow$ update by
corresponding h_t incrementing by α_t

- Functional Gradient Descent

function of function

$$L(F) = L(F(x_1), \dots, F(x_m)) = \sum_i \exp(-y_i F(x_i))$$

↳ gradient decent: $F \leftarrow F - \alpha \nabla_F L(F)$
restricted: $F \leftarrow F + \alpha h_t \quad \Big| \Rightarrow h_t$ close to $-\nabla_F L(F)$

↳ Equivalent to Ada Boost

- Estimate Prob ($y_{+1} | x$)

↳ empirical version → x, y random from true dist

$$\mathbb{E}_{x,y} [e^{-yF(x)}] = \mathbb{E}_x \left[\Pr(y_{+1} | x) e^{-F(x)} + \Pr(y_{-1} | x) e^{F(x)} \right]$$

$$F(x) \underset{\min}{=} \frac{1}{2} - \ln \left(\frac{\Pr(y_{+1} | x)}{\Pr(y_{-1} | x)} \right) \Rightarrow \Pr(y_{+1} | x) = \frac{1}{1 + e^{-2F(x)}}$$

- Calibration Curve



- Regularization

$$L(\lambda) = \sum \exp\left(-y_i \sum_j \lambda_j g_j(x_i)\right)$$

$$l_1\text{-reg: } \min L(\lambda) \text{ s.t. } \|\lambda\|_1 \leq B$$

$$\equiv \min L(\lambda) + B \|\lambda\|_1$$

↳ if $\alpha \rightarrow 0$: results looks like Ada Boost

* Information - Geometric View

- Loss Minimization \rightarrow Functions computed on ada boost \rightarrow weights on weak classifier
 \Downarrow dual
Info-Geometric \rightarrow distributions \rightarrow weights on examples
- Iterative - Projection Alg
 - Find point closest to x_0 in set $\rightarrow P = \{\text{intersection of hyperplane}\}$
 - pick $x_0 \rightarrow$ project on hyperplane \rightarrow converge if $P \neq \emptyset$
 - Ada Boost is Iterative - Projection
 - \hookrightarrow points = dists over training examples
 - \hookrightarrow distance = relative entropy \rightarrow distance between two distributions
 - \hookrightarrow reference point $x_0 =$ uniform dist
 - \hookrightarrow hyperplanes defined by all possible weak classifiers \rightarrow looking for hardest dist.
 - Alg.
start from uni. dist \rightarrow pick weak classifier \rightarrow entropy projection on hyperplane

* Multi-class Problem

- Direct Approach $\rightarrow H_{\text{final}} = \arg \max_{y \in \{1, \dots, k\}} \sum_t \kappa_t$
Ada Boost.M1 \rightarrow t: $h_t(x) = y$
 \searrow bad result for "weak" weak classifier
- One-Against-All \rightarrow k - Binary class
Ada Boost.MH
- Output Coding

* Ranking Problem

- Learn to Rank \Rightarrow Rank Boost \Rightarrow lots of decision of A better than B
from example

* Confidence Rated Prediction

- Ideal prediction \rightarrow one side of prediction is always + example
 \searrow the other side is not important
- Hard prediction \rightarrow bad predictions \rightarrow need clean up
 \searrow speed \downarrow
- introduce confidence $\rightarrow |h_t(x)| =$ confidence
 $\searrow \text{sgn}(h_t(x)) =$ prediction

* Optimal Accuracy

- AdaBoost can converge to Bayes optimal
 - enough data, run for many runs, good weak classifiers
 - Universally consistent
- Noise
 - there are data sources that AdaBoost fails miserably
 - but anyway is better than random \rightarrow err ≤ 50
 - works bad \rightarrow not rich enough classifiers
 - \rightarrow prone to noise \rightarrow solution: branching program
 - \rightarrow regularization may not help

* Optimal Efficiency

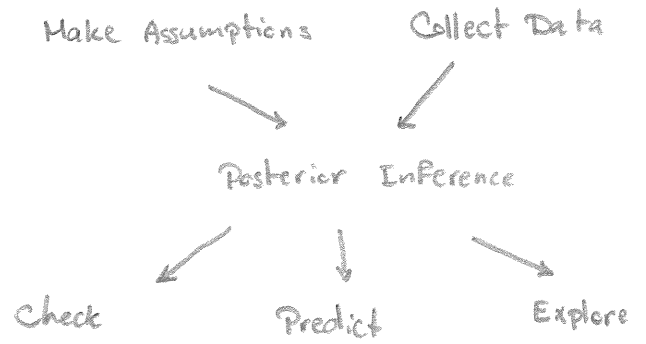
- AdaBoost works no better than boost-by-majority
 - AdaBoost train error is like Chernoff approx of BSM

Topics Models

(Elei - Princeton)

* Topic Models \rightarrow Bayesian Model

- Specifying models
- Model Selection
- Approximate posterior inference
- Using and Evaluating Models



□ Bayesianly \rightarrow word appear in [Robin, 1984]

* Latent Dirichlet Allocation (LDA)

• Intuition: Docs exhibit multiple topics

- Topic: Dist. over words

- Doc: New dist over topics for each doc. \rightarrow Bag-of-words assumption

- Topics are fixed over all docs. But each doc may have different topics

\rightarrow Mixture

\rightarrow Topics fixed for all docs, docs have different dist. of topic

- Simplex

- Doc \rightarrow Empty Topic \rightarrow According to dist of topics \rightarrow Pick a topic

\rightarrow According to topic dist \rightarrow pick a word $(z_{d,n})$ \rightarrow Repeat it for all words in topic

• A graphical models

- edges \rightarrow possible dependence

- nodes \rightarrow random variable (shaded = observed)

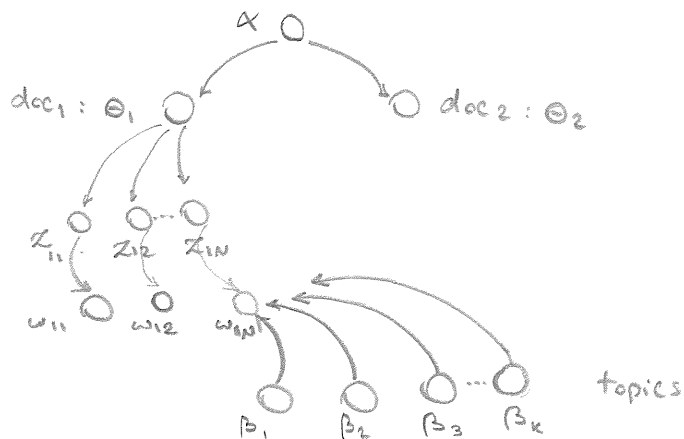
- plane (\square) \rightarrow hyper parameters

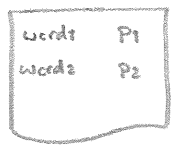
- Encodes dependence assumption

- $z_{d,n}$ \rightarrow depends on topics and

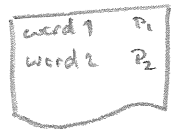
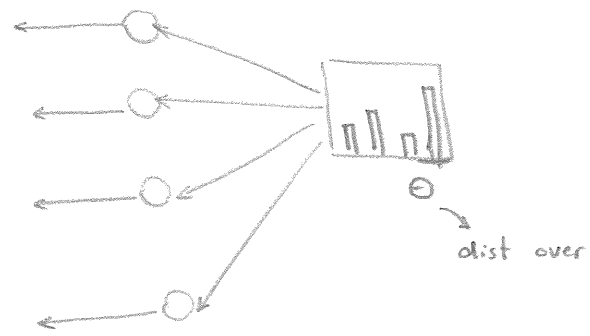
\rightarrow index of topic

\rightarrow drawn on θ





word in doc



$$P(\beta, \theta, z | w) = \frac{P(\beta, \theta, z, w)}{P(w)} = \frac{P(\beta, \theta, z, w)}{\int_{\beta, \theta, z} P(\beta, \theta, z, w)}$$

$$P(\beta, \theta, z, w) = \left(\prod P(\beta_i | \eta_i) \right) \left(\prod_{d=1}^D P(\theta_d | \alpha) \prod_{n=1}^N P(z_{d,n} | \theta_d) P(w_{d,n} | \beta_{1:k}, z_{d,n}) \right)$$

✓ jstor → scan docs → OCR them → index them

- conjugate → if prior comes from a family of dist. (Dirichlet for example) → sampled posterior is from the same family

• Dirichlet Distribution

$$P(\theta | \alpha) = \frac{\Gamma(\sum \alpha_i)}{\prod \Gamma(\alpha_i)} \prod \theta_i^{\alpha_i - 1}$$

$$E[\theta] = \frac{\alpha}{\sum \alpha_i} \rightarrow \checkmark \text{Dir}(\alpha, \alpha, \dots, \alpha) \Big|_{\alpha=1}^{10} \Rightarrow E[\theta] = \frac{1}{10}$$

$\alpha \uparrow \Rightarrow$ closer to the mean

$\alpha \downarrow \Rightarrow$ more sparsity $\Rightarrow \sigma^2 \uparrow$

$\alpha \rightarrow 0 \Rightarrow$ choose on prior and give it all the mass (=1)

→ ✓ Movie: powers of ten

• Why LDA works?

- For each doc, allocate words to as few topics as possible.

$$P(\beta, \theta, z | w) \propto P(\beta, \theta, z, w)$$

$$= P(\beta) \prod_d P(\theta_d) \left(\prod_n P(z_{d,n} | \theta) P(w_{d,n} | z_{d,n}, \beta) \right)$$

$P(z_{d,n} | \theta) \downarrow \Rightarrow P(\beta, \theta, z | w) \downarrow$

more topics $\Rightarrow P(z_{d,n} | \theta) \downarrow$

← it sits on simplex

- For each topic, assign high probability to as few terms as possible

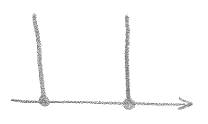
↳ same reason...

- if we have only 1 topic \Rightarrow high prob \rightarrow more topic \Rightarrow posterior \downarrow

if we have only 5 topic \Rightarrow high word prob \rightarrow more words \Rightarrow posterior \downarrow

* Extending LDA: Correlated and Dynamic Topic Model

• LDA \rightarrow hidden assumption \rightarrow component are nearly independent



• Logistic Model \rightarrow capture dependence between component
 \rightarrow is not conjugate

• topics in the course of time are dependant $\beta_{t,k} | \beta_{t-1,k} \sim \mathcal{N}(\beta_{t-1,k}, 1/\alpha^2)$

• Dynamic Topic Model:

- if some new topic introduces because of low variance estimation before it was penalized in previous model \rightarrow so we drift the mean toward that a little bit to explain that

- so we need to ease the acceptance of new topics in future, so they should give feedback to previous stages.

* Bayesian Non parametric Models

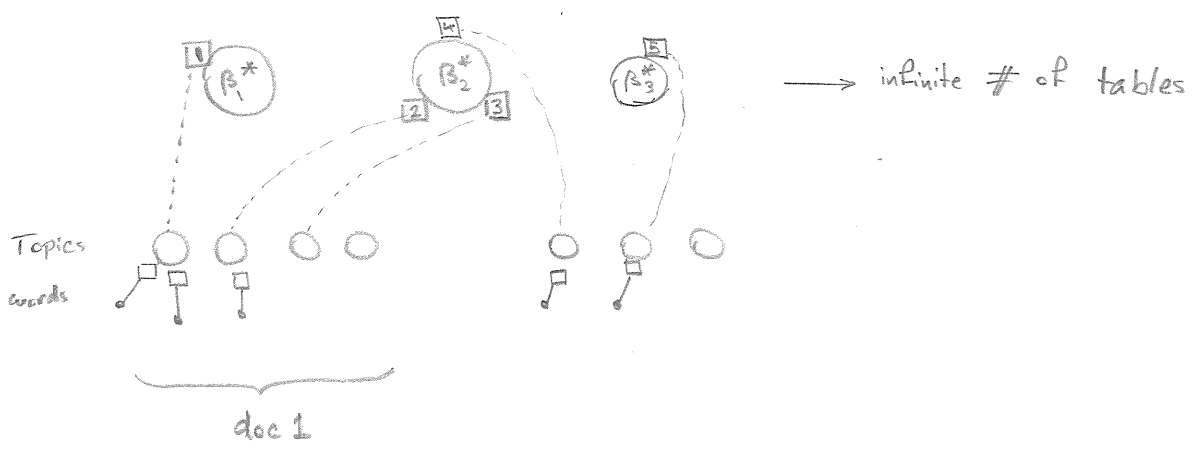
• Topic Models \rightarrow fixed number
 \rightarrow regularization parameter \rightarrow CR Model selection

\rightarrow using non-parametric models

• Non Param. Models (BNPM) \rightarrow extract new topics in corpus in future time.

• Chinese Restaurant Process (CRP)

• The chinese restaurant franchise



Graphical Model :

(Wainwright - Berkeley)

* Computational Challenges :

- Efficiently compute MAP

$$\hat{x} = \arg \max_{x_1, \dots, x_N} \underbrace{P(x_1, x_2, \dots, x_N | y_1, \dots, y_N)}_{\text{observation}}$$

Shannon Experiment of Generating English Sentences with Markov Networks

- Normalization Constants

$$Z = \sum_{x \in X^N} \prod_{c \in C} \psi_c(x_c)$$

* Efficiently Solving MAPs \rightarrow Belief propagation, message passing, max-product

- No cycles in graphs
- Compute MAP on tree \rightarrow Extension on Viterbi Algorithm

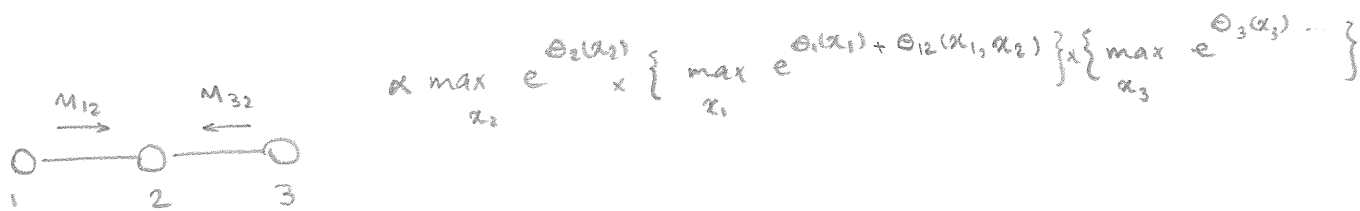
$$\hat{x} = \arg \max_{x \in X^N}$$

- What info each node need to receive from neighbor to do the computation

\rightarrow Message: vectors of number, blue arrows

\rightarrow Divide & Conquer principle

$$\max_{x_1, x_2, x_3} P(x_1, x_2, x_3) \propto \max_{x_1, x_2, x_3} e^{\theta_1(x_1)} e^{\theta_2(x_2)} e^{\theta_3(x_3)} e^{\theta_{12}(x_1, x_2)} e^{\theta_{23}(x_2, x_3)}$$



brute force: (2^n) \rightarrow message passing: (n^2)

won't work if there is an edge from 1 to 3 \rightarrow cycle \rightarrow cannot simplify

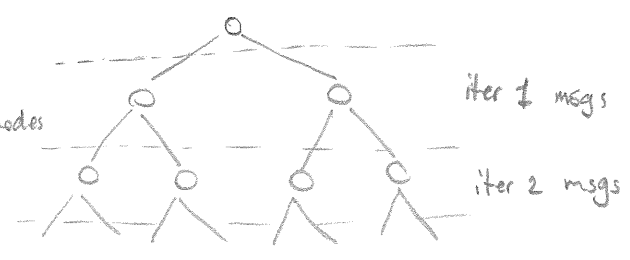
- Message start from leaf \rightarrow start to complete subtrees \rightarrow when a node has enough info
- iter # \rightarrow diameter of graph: longest path
- parallel version of Max-Product vs. flooding schedule: the trends in computation and implementation: 2 sweep
- Generalization \rightarrow Junction Trees

* Max-product on graph with cycles

- Messages as rumors → inexact best guesses
- ⊕ when you hear a rumor → From more people: stronger belief
 - ↘ From a far distance: it gets weaker → long chains
- No general guarantee → Can algorithm lie? yes! → converge to bad answers

- Computation tree

↓
there is repetitive nodes



whereas the information flowing in graph through iterations

- Solving Max-product on the computation tree

- ↳ same potential graph
- ↳ lots of copies on one nodes → some visited more often → higher degree
 - ↳ not balanced in computation
 - ↳ having higher weights

- Fix repetitive node problems → instead of telling a lie says I'm confused.

↳ weighting the edges

* Tree-reweighted Max-Product

- I won't get back any of my own messages.
- Visit every node and edge the same # of times → by weighting them
- Where the weights come from? → one possibilities

↳ put all spanning tree with same prob. in a hat!

↳ what is the prob. that an edge appear in spanning trees → count 'em!

- converge: if you schedule messages correctly → won't lie → MAP result guaranteed or say i'm confuse

• rewrite the MAP as Linear Programming

$$\hat{x} = \arg \max_{x_1, \dots, x_N} \left\{ \sum_{S \in V} \theta_S(x_S) + \sum_{(s,t) \in E} \theta_{st}(x_s, x_t) \right\}$$

$J(x, \theta) \rightarrow$ cost function

instead of max over x , maximize on distribution supporting them

$$\hat{x} = \max_{q(\cdot)} \left\{ \sum q(x) J(x, \theta) \right\} \quad \begin{matrix} q(\cdot) \geq 0 \\ \sum q = 1 \end{matrix}$$

$$\hat{x} = \max \left\{ \sum_{s \in Y} \sum_{x_s} \mu_s(x_s) \Theta_s(x_s) + \sum_{(s,t) \in E} \sum_{x_s, x_t} \mu_{st}(x_s, x_t) \Theta_{st}(x_s, x_t) \right\}$$

local marginals
 $\mu_s(x_s), \dots, \mu_{st}$

$$\begin{bmatrix} \mu_s(x_s) \\ \vdots \\ \mu_s(x_{m-1}) \end{bmatrix} \geq 0$$

$$\sum \mu_s = 1$$

$$\begin{bmatrix} \mu_{st}(x_s, x_t) & \dots \\ \vdots & \ddots \\ \mu_{st}(x_{m-1}, x_{m-1}) \end{bmatrix}$$

$$\sum_{x_t} \mu_{st}(x_s, x_t) = \mu_s(x_s)$$

↳ constraints of linear programming

↳ messages in max-products \equiv Lagrange product in LP

• Marginal polytope

$$M(G) = \left\{ \mu_s, \mu_{st} \mid \begin{array}{l} \mu_s(x) = \sum_{x_u, u \neq s} q(x_1, \dots, x_N) \\ \mu_{st}(x) = \sum_{x_u, u \neq t, s} q(x_1, \dots, x_N) \end{array} \text{ for some } q(\cdot) \right\} \rightarrow \text{convex set}$$

• pseudo marginals

$$L(G) = \left\{ T_s, T_{st} \mid T_s(x_s) \geq 0, \sum_{x_s} T_s(x_s) = 1, \sum_{x_t} T_{st}(x_s, x_t) = T_s(x_s) \right\}$$

semi-definite relaxation

• How many of polytope halfspaces we need? $L(G)$ has few constraint \rightarrow easy to optimize
 How many constraint needed in $M(G)$?

□ $T \notin M(G)$ by semi definite constraints

$$x_s \in \{0, 1\} \quad E(x_s) \stackrel{?}{=} 0.5 \text{ for } s: 1, 2, 3$$

and figure in P27/35

$$s = 1, 2, 3 \quad E(x_s, x_t) \stackrel{?}{=} \begin{cases} 0.9 & \text{for } (s,t) \in \{(1,2), (2,3)\} \\ 0.1 & \text{for } (s,t) = (1,3) \end{cases}$$

$$E \left[\begin{pmatrix} 1 \\ x_1 \\ x_2 \\ x_3 \end{pmatrix} (1, x_1, x_2, x_3) \right] = E \begin{bmatrix} 1 & x_1 & x_2 & x_3 \\ x_1 & x_1 & x_1 x_2 & x_1 x_3 \\ x_2 & x_1 x_2 & x_2 & \dots \\ x_3 & x_1 x_3 & x_2 x_3 & x_3 \end{bmatrix}$$

binary values = $x_i^2 = x_i$

$$= \begin{bmatrix} 1 & 0.5 & 0.5 & 0.5 \\ 0.5 & 0.5 & 0.4 & 0.4 \\ 0.5 & 0.4 & 0.5 & 0.1 \\ 0.5 & 0.4 & 0.1 & 0.5 \end{bmatrix}$$

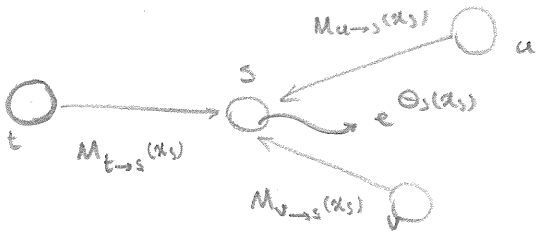
\rightarrow eigen value should be positive because the first E is rank one \rightarrow but it is not P.S.D!

$$\stackrel{?}{\geq} 0$$

↳ positive semi-definite

* Sum-product message passing

• replacing max operator by sum in all formula



$$P(x_s = \bar{x}_s) \propto \exp\{\Theta_s(\bar{x}_s)\} \prod_{u \in N(s)} M_{u \to s}(\bar{x}_s) \quad \text{P.4}$$

Kalman Filters \rightarrow special case with Gaussian Dist

\hookrightarrow Msgs: means + covariances

\hookrightarrow use these relations and put Gaussian in them too.

Fast Fourier Transform \rightarrow special case too

[Aji & McEliece, 2007 \rightarrow survey]

using bethe entropy approx.

- information on edges
- using entropy as cost function

Entropy and Markov Chain



$$P(x_1, x_2, x_3) = \mu_1(x_1) \mu_{2|1}(x_2|x_1) \mu_{3|2}(x_3|x_2)$$

$$= \mu_1(x_1) \times \frac{\mu_{12}(x_1, x_2)}{\mu_1(x_1)} \times \frac{\mu_{23}(x_2, x_3)}{\mu_2(x_2)}$$

$$= \mu_1(x_1) \times \mu_2(x_2) \times \mu_3(x_3) \times \frac{\mu_{12}}{\mu_1 \times \mu_2} \times \frac{\mu_{23}}{\mu_2 \times \mu_3}$$

$$H(P) = - \sum_{x_1, x_2, x_3} P(x_1, x_2, x_3) \log P(x_1, x_2, x_3)$$

$$H_{\text{Bethe}}(\mu, P) = \sum_{s \in V} H_s(x_s) = \sum_{(s,t) \in E} P_{st} I_{st}(\mu_{st})$$

\hookrightarrow mutual info

why entropy arise in duality

$$P(x_1, \dots, x_N; \theta) = \frac{1}{Z(\theta)} \exp\left\{ \sum_{c \in C} \theta_c(x_c) \right\} \quad \rightarrow \text{exponential family}$$

$\Phi(x) = x \rightarrow E(x) = \hat{\mu}$ = sample mean

Max Entropy & Max Likelihood

$$\alpha \in \{0, 1\} \quad P(\alpha, \theta) = \exp\{\theta \alpha - A(\theta)\} \quad A(\theta) = \log Z(\theta) = \log(1 + e^\theta)$$

$$\mu = P(\alpha = 1) \quad -H(\mu) = \mu \log \mu + (1-\mu) \log(1-\mu) = A^*(\mu)$$

$(A, A^*) \rightarrow$ conjugate dual pair

$$A^*(\mu) = \sup_{\theta \in \mathbb{R}} \{\theta \mu - A(\theta)\}$$

derivative w.r.t.

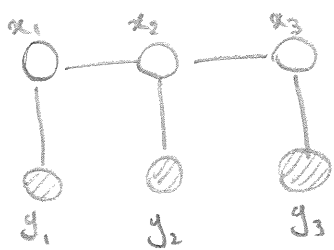
$$\mu - \frac{e^\theta}{1+e^\theta} = 0$$

} Moment Matching problem
 $1+e^\theta = P(\alpha=1; \theta)$

if $\mu \in (0, 1)$ has an answer: $A = \log \frac{\mu}{1-\mu}$

$$\Rightarrow A^*(\mu) = \begin{cases} \mu \log \mu + (1-\mu) \log (1-\mu) & \text{if } 0 < \mu < 1 \\ +\infty & \text{otherwise} \end{cases}$$

• why do we care about partition function?



$$P(\alpha, \theta) \propto \prod_{c \in C} e^{\theta_c(\alpha_c)} \quad \text{prior}$$

$$P(y|\alpha) \propto \prod_i P(y_i|\alpha_i)$$

$$P(\alpha|y, \theta) = \frac{P(\alpha; \theta) P(y|\alpha)}{P(y, \theta)}$$

$$H_{\text{Bethe}}(\mathcal{T}, P) = \underbrace{\sum_{s \in V} H_s(\mu_s)}_{\text{singleton entropy}} - \sum_{(s,t) \in E} \underbrace{P_{st}}_{\text{edge weight}} \underbrace{I_{st}(\mu_{st})}_{\text{mutual info}}$$

• we need upper bound

$$A(\theta) = \log Z(\theta) = \log \sum_{\alpha} \prod_c e^{\theta_c(\alpha_c)}$$

$$A\left(\sum_T f(T) \theta(T)\right) \leq \sum_T f(T) A(\theta(T)) \rightarrow \underline{A} \text{ is convex, apply Jensen's}$$

- many trees in a graph \rightarrow number grows exponentially

- use duality \rightarrow put exponential explosion aside

* Learning Graphs from Data for pairwise models

$$Q(\alpha_1, \dots, \alpha_p; \theta) = \frac{1}{Z(\theta)} \exp \left\{ \sum_{s \in V} \theta_s \alpha_s^2 + \sum_{(s,t) \in E} \theta_{st} \alpha_s \alpha_t \right\}$$

(a) $\alpha_s \in \mathbb{R} \rightarrow$ multivariate Gaussian (zero-mean)

(b) $\alpha_s \in \{-1, +1\} \rightarrow$ Ising Model

Unknown: edge structure and weights

Parameter Matrix: $\Theta \in \mathbb{R}^{p \times p} \rightarrow p = \#$ of nodes

Learning Problem: collect sample $X_{i \cdot} = (X_{i1}, \dots, X_{ip})$

$$X_i^n = \begin{bmatrix} \alpha_i^T \\ \vdots \\ \alpha_n^T \end{bmatrix} \in \mathbb{R}^{n \times p} \quad (\text{or } \in \{0, 1\}^{n \times p})$$

Wanted: Support on $\Theta \rightarrow$ non-zeros } all edges
no false edges

$$\text{supp}(\Theta) = \{(s,t) \mid \Theta_{(s,t)} \neq 0\}$$

$$\text{Estimator: } X_i^n \mapsto \hat{\Theta}$$

US Senate Net $\rightarrow p$: # of senators
 $\rightarrow n$: # of votes in DB

• with iid sample

$$e(\theta; X^n) = -\frac{1}{n} \sum \log Q(\cdot) = \log Z(\theta) - \sum_s \hat{\mu}_s \theta_s - \sum_{(s,t)} \hat{\mu}_{st} \theta_{st}$$

from data \leftarrow empirical count

- depends on Z : partition function
- for gaussian case \rightarrow covariance selection: log-determinant program
- generally use pseudo-likelihood \rightarrow fake likelihood
- for tree: problem of ML \Rightarrow Max Spanning Tree

Markov Chain:

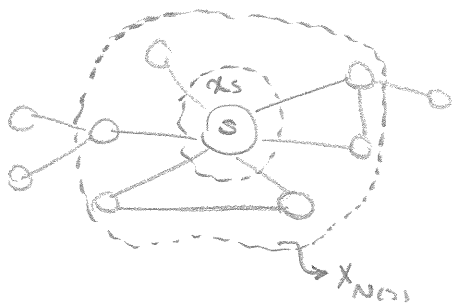


$$X_{\text{Past}} \perp\!\!\!\perp X_{\text{Future}} \mid X_3$$

if we cut $x_3 \rightarrow$ the graph breaks \rightarrow cutset

if we remove cutset \rightarrow the two partitions should be independent

• what kind of cutset would be useful?



- reduce to neighborhood selection

$$N(s) = \{t \in V \mid (s,t) \in E\}$$

$N(s) \rightarrow$ is a cut set

$$X_{N(s)} = \{x_t \mid t \in N(s)\} \rightarrow \text{Markov Blanket}$$

• Finding Cutsets

- if we know the size \rightarrow Naive way = test all sets of nodes, cut them and check computational independence
- advanced papers

• Graph Selection via neigh. regression

$$\hat{\Theta}[s] = \underset{\Theta}{\text{argmin}} \left\{ -\frac{1}{n} \sum \text{local log likelihood} + \lambda \times \text{regularization} \right\}$$

$$l_1\text{-norm} = \|\Theta\|_1 = \sum_{j=1}^p |\theta_j| \rightarrow \text{induce sparsity} \rightarrow \begin{cases} \text{few large numbers} \\ \text{lots of zeros} \end{cases}$$

it could be logistic regression

binary model
conditional likelihood

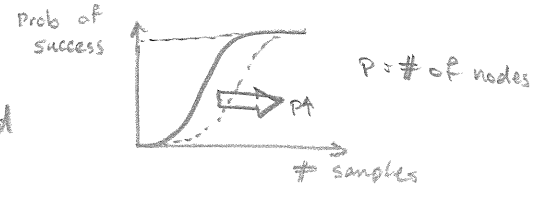
$$\Theta(x_s=1 | x_{N(s)}) = \frac{\exp(\sum_{t \in N(s)} \Theta_{st} x_t)}{1 + \exp(\tilde{\mu})}$$

Markov property

$$x_{N(s)} = x_{N(s)}$$

- bigger graphs need more samples \rightarrow harder to estimate more edges

star graph \rightarrow show method works
 \rightarrow how many samples do we need



degree $\uparrow \Rightarrow$ more samples needed

d = maximum degree, p = graph size, n = # samples

For guarantee \rightarrow no false positive
 \rightarrow no false negative (no unfound edge)

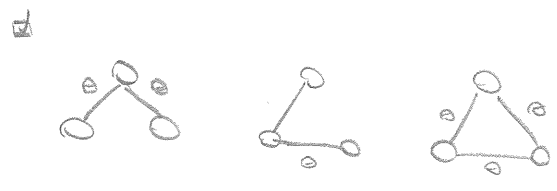
$$\frac{n}{d^3 \log P} > c_1$$

$$\lambda n > c_1 \sqrt{\frac{\log P}{n}}$$

with prob greater than $1 - 2e^{-c_2 \lambda^2 n}$

• Challenges of Graph Selection

- Lower bound on edge weights \Rightarrow SNR big enough
- Need upper bound on $|\Theta_{st}|$



$(x_1, x_2, x_3) \in \{-1, +1\}^3$
 $n \rightarrow +\infty$
these three graph will be similar because of large weights

• Graph selection as channel coding

- environment give you the info via noisy samples
- $\mathcal{G}_{p,d}$ = all graphs on p vertices with max-degree channel

$$\Theta_{min} = \min_{(s,t) \in E} |\Theta_{st}|$$

$$w(\Theta) = \max_{S \in V} \sum_{t \in N(s)} |\Theta_{st}|$$

graph size $\uparrow \Rightarrow$ # of samples needed $\uparrow \uparrow \uparrow$ exponentially

$$n < \max \left\{ \frac{\log P}{2\Theta_{min} \tanh(\Theta_{min})}, \dots \right\} \Rightarrow w(\Theta) \geq d\Theta_{min} \Rightarrow \Theta_{min} \leq \frac{1}{d}$$

$$n \gtrsim d^2 \log P \leftarrow \approx \frac{\log P}{\Theta_{min}^2} \rightarrow \text{no really small } \Theta_{min} \leftrightarrow \Theta_{min} \approx \sqrt{\frac{\log P}{n}}$$

* Problems in High Dimension

- Optimization
- Integration
- Sampling
- Rounding
- Learning

* High D Problem

- Hard \rightarrow intractable for \rightarrow linear func.
- \searrow polytope & quadratic func.

- NP \rightarrow is there a solution or not?
- #P \rightarrow what is the number of solution

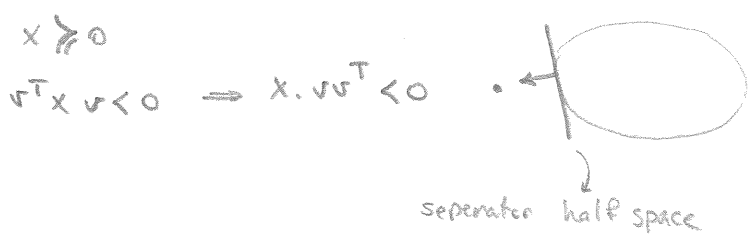
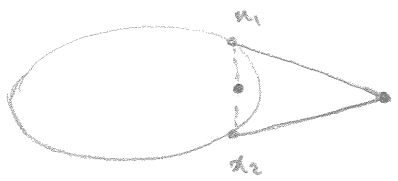
- Ellipsoid Algorithm \rightarrow Optimization for convex sets
- DFK Algorithm \rightarrow Integration

• How to specify a convex set?

- explicit list of constraints
- set of positive semi definite matrices?

- a point that is not in the set \rightarrow there is a hyper plane that separate those \rightarrow
 there is a unique point that is closest to the point

برهان صحت نتیجه: اگر x_1 و x_2 هم در مجموعه باشند، $\frac{x_1+x_2}{2}$ نیز در مجموعه است.



• Majority Algorithm

- input $x_1, \dots, x_k \rightarrow x_{k+1}$ what should we decide?
 $f(x_1), \dots, f(x_k)$?

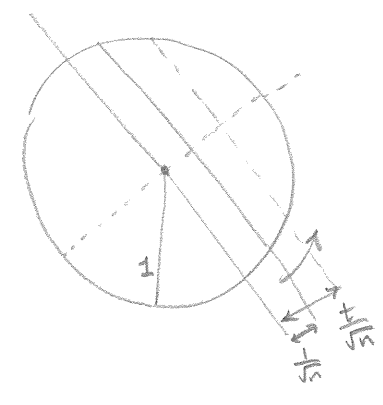
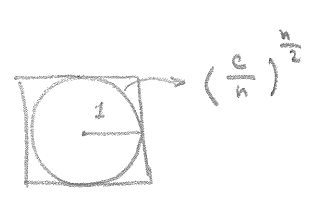
- majority of previous one? x
- nearest neighbor alg \checkmark

- set a hyper plane $S_k = \{w \mid \forall i, w^T (\text{sgn}(x_i) x_i) > 0\}$ \rightarrow err $\leq b_n \checkmark$
 \rightarrow set is large x

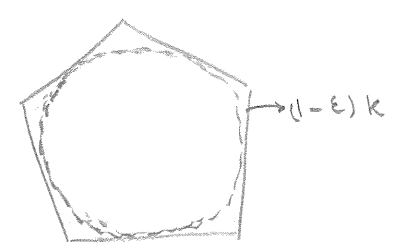
• Random Algorithm → Choose a point and set hyperplane accordingly

→ $E(\# \text{ wrong guesses}) \leq b_n$ # of bits of sample

• Volume of sampling area



$\text{Vol}(\text{outside } \frac{t}{5}) \leq 2e^{-\frac{t^2}{2}}$

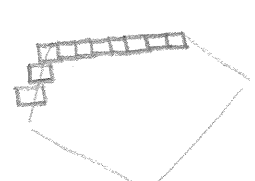


• Brunn Minkowski Inequality

$$\left(\prod (a_i + b_i) \right)^{\frac{1}{n}} \geq \left(\prod a_i \right)^{\frac{1}{n}} + \left(\prod b_i \right)^{\frac{1}{n}}$$

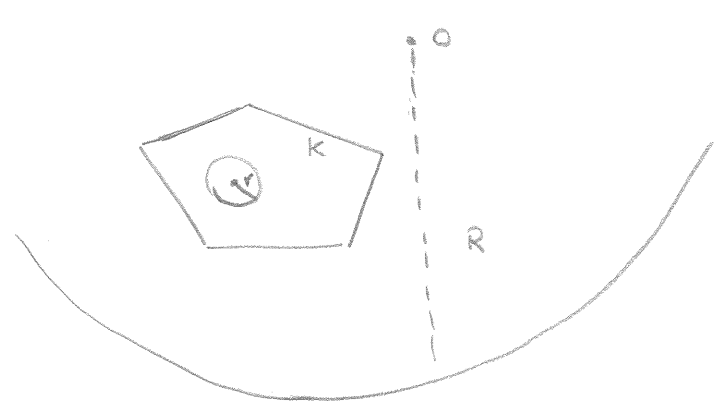
$$1 \geq \left(\prod \frac{a_i}{a_i + b_i} \right)^{\frac{1}{n}} + \left(\prod \frac{b_i}{a_i + b_i} \right)^{\frac{1}{n}}$$

because: $\frac{1}{n} \sum x_i \geq \left(\prod x_i \right)^{\frac{1}{n}}$
 other proof by taking log



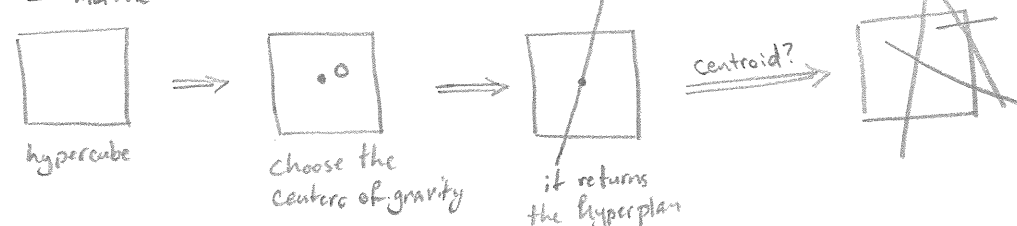
• Convex Feasibility

K is isotropic
 if $E_K(X) = 0$
 $E_K(XX^T) = I$



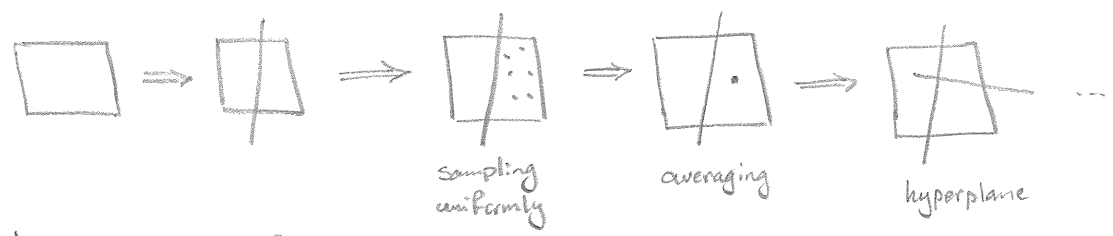
• How to choose oracle query

- naive



start with R^n and ends with r^n → how many steps? $\log\left(\frac{R}{r}\right)^n = n \log \frac{R}{r}$
 we don't know how to calculate centroids ... $\neq P$ solutions

- Bertsimas



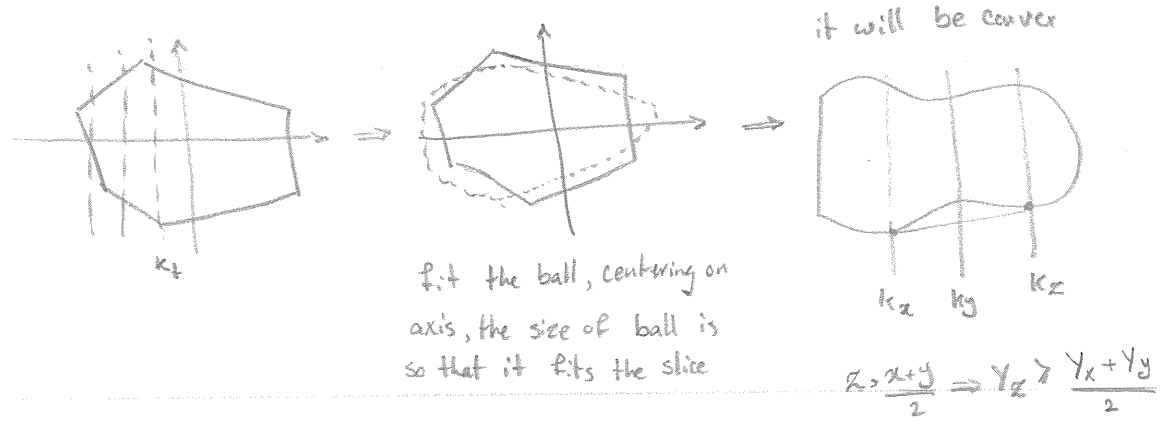
how many samples? are we doing any progress?

- Centroid Alg

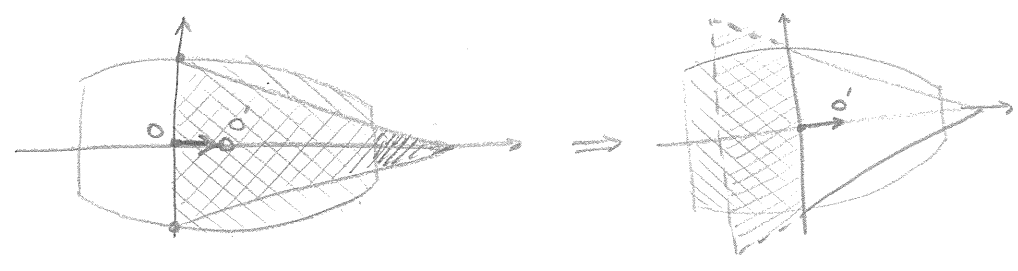
↳ cut vol by a constant fraction (at least $\frac{1}{e}$) $\Rightarrow \text{vol}(K_{i+1}) \geq \frac{1}{e} \text{vol}(K_i)$

↳ number of steps $\log_{\frac{1}{e}} \left(\frac{R}{r}\right)^n = O(n \log \frac{R}{r})$

↳ centroid cuts are balanced



$$\text{vol}(K_z)^{\frac{1}{n-1}} \geq \frac{1}{2} \text{vol}(K_x)^{\frac{1}{n-1}} + \frac{1}{2} \text{vol}(K_y)^{\frac{1}{n-1}} \xrightarrow{f(n-1) r_z^{n-1}} y_z \geq \frac{1}{2} y_x + \frac{1}{2} y_y$$



make the triangle in away that the area is equivalent to the shape \Rightarrow shift center of gravity right

extend lines so the left part equals the tails of triangle \Rightarrow again shift the center right

↳ how many sample?

* Volume Estimation

• it is told that we have a unit ball, a convex set is contains in that, you should guess the volume.

all you can do is you can ask if a point is in the ball or not.

the convex hull of points is the answer \rightarrow number of questions required is required

• volume from sampling

$$k_0 = \text{ball}$$

$$\vdots$$

$$k_m = k$$

$$\frac{\text{vol}(k_{i+1})}{\text{vol}(k_i)} \sim \text{small}$$

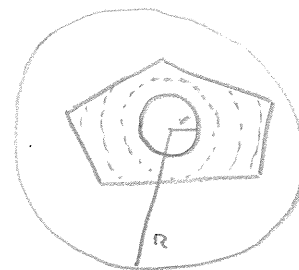
$$\text{claim } \text{vol}(k_{i+1}) \leq 2 \text{vol}(k_i)$$

$$\text{vol}(k) = \text{vol}(B) \times \frac{\text{vol}(k_1)}{\text{vol}(k_0)} \times \dots \times \frac{\text{vol}(k_m)}{\text{vol}(k_{m-1})}$$

$$k_{i+1} = k_i \cap 2^{\frac{i+1}{n}} B \subseteq 2^{\frac{i}{n}} (k_i \cap 2^{\frac{i}{n}} B)$$

$$\begin{aligned} \text{vol}(k_{i+1}) &\leq \text{vol}(2^{\frac{i}{n}} (k_i \cap 2^{\frac{i}{n}} B)) \\ &= \text{vol}(2^{\frac{i}{n}} k_i) = 2 \text{vol}(k_i) \end{aligned}$$

$$\text{total \# of sample} = m \times \frac{m}{\epsilon^2} = O^*(n^2)$$



* Problem

- multiple prob. dist
- solving \rightarrow no need to individual prob

$$r(x) = \frac{p(x)}{q(x)}$$

• Vapnik's principle: should not solve a more general problem as intermediate step

• Kernel least-squares (KLS) needed only!

- importance sampling
- KL divergence estimation
- Mutual information
- Conditional prob. estimation

* Density Ratio Estimation Method:

• problem: estimate r from two set of iid samples

• Naive

1. density estimation
2. ratios of densities

\rightarrow separated steps \times

• Probabilistic Classification

- separate numerators and denominator samples by a probabilistic classifier
- using logistic regression \rightarrow model correct \Rightarrow achieve min
- not reliable for misspecified models
- parallel processing

• Moment Matching

- Match moments of \hat{P}_{nu} and $P_{nu} \rightarrow \hat{P}_{nu} = \hat{r}(x) P_{de}(x)$

\rightarrow better to do in higher dimensions

\rightarrow computation \uparrow

\rightarrow need kernel trick \rightarrow work well if σ of gaussian kernel is good

- In practice: using QP for optimization

- Setting σ (\odot) \rightarrow median distance between samples $\checkmark \rightarrow$ fails in multimodal

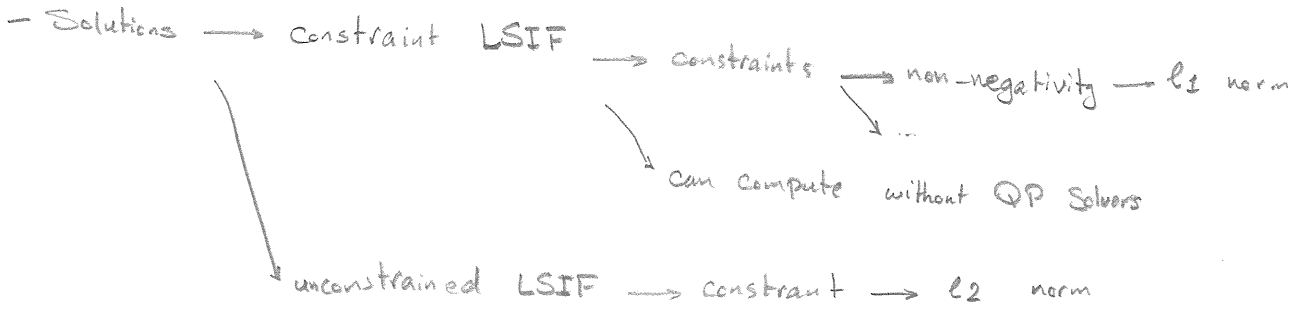
• Density Fitting:

- KL Importance → Minimize KL divergence from $P_{true}(x)$
- Gradient Ascent on KL → project onto feasible region

• Density Ratio Fitting:

- Minimize Squared-loss

↳ like regression, but different cause r_n should be sampled



* Usage:

- learning under Covariate shift → train / test dist. different
- ↳ target function fix

- Covariance shift → $n \rightarrow \infty$ can use ordinary least sq.
- ↳ $n \downarrow$ not consistent

↳ use importance weighting → importance weighing Least Squares

↳ bias-var problem → flattened importance weight

- Model Selection → Akaike info. criterion

↳ Subspace info. "

↳ Cross validation

* Data

- For sciences & non-scientific applications (e.g. data mining)
- data drives the revolutions → complex dynamical systems
 - purely data driven things
 - ✓ Market place
 - ✓ Mining DBs for information
- broad range of ways to tackle data
 - signal processing
 - statistical methods → time and frequency analysis
 - dimensionality reduction & PCA
 - machine learning
 - integrating dynamics
 - data collection
 - teach computer how to see bigger trends in data
 - construct models to predict future states

* Overview of Course

- time frequency analysis
 - noisy signals in time
 - representing them in time & frequency
 - prediction
- machine learning
 - computer science + statistics → pattern recognition
 - data features
 - singular value decomposition
- complex dynamical systems
 - dimension reduction
 - equation free modeling
 - dynamic mode decomposition

* Time - Frequency Analysis

↳ Fourier Transform

↳ Wavelets

• Fourier Transform

- idea: I can take a given function and I can represent any function as a sum of cosine and sines.

$$f(x) = \frac{a_0}{2} + \sum_{n=1}^{\infty} (a_n \cos nx + b_n \sin nx) \quad x \in [-\pi, \pi]$$

- how to compute these a_i and b_i ?

↳ multiply both sides to $\cos mx$ → then integrate it from $-\pi$ to π

$$\int_{-\pi}^{\pi} f(x) \cdot \cos mx \cdot dx = \frac{a_0}{2} \int_{-\pi}^{\pi} \cos mx \cdot dx$$

$$+ \sum_{n=1}^{\infty} a_n \int_{-\pi}^{\pi} \underbrace{\cos nx \cos mx \cdot dx}_{\substack{\text{if } n \neq m \rightarrow 0 \\ \text{if } n = m \rightarrow \pi}} + b_n \int_{-\pi}^{\pi} \underbrace{\sin nx \cos mx \cdot dx}_{\substack{n, m \in \mathbb{Z} \\ \sin nx \cdot \cos mx = 0}} \cdot dx$$

orthogonal function

↳ orthogonality in function space → $\cos 1x, \cos 2x, \dots$ are directions in func. space

orthogonal to $\sin 1x, \sin 2x, \sin 3x, \dots$

$$\rightarrow a_n = \frac{1}{\pi} \int_{-\pi}^{\pi} f(x) \cdot \cos nx \cdot dx$$

$$\rightarrow b_n = \frac{1}{\pi} \int_{-\pi}^{\pi} f(x) \cdot \sin nx \cdot dx$$

- complex representation

$$f(x) = \sum_{-\infty}^{\infty} c_n e^{in \frac{\pi}{L} x} \quad x \in [-L, L]$$

$$c_n = \frac{1}{2L} \int_{-\pi}^{\pi} f(x) \cdot e^{-in \frac{\pi}{L} x} \cdot dx$$

- Fourier says that we can represent a function as sines and cosines

↳ we think of $f(x)$ as a time signal.

↳ we care about coefficients → $n=1: c_1, n=2: c_2 \dots$

↳ $n=1 \rightarrow a_1 \rightarrow \cos x \rightarrow$ has a frequency

↳ take a signal in time and show it as a collection of frequencies components

- MATLAB

we have a function u defined (as vector)

transform to Fourier domain $ut = \text{fft}(u) \rightsquigarrow$ fast Fourier transform

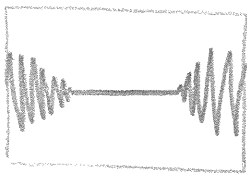
fft \rightarrow late 60's by Tockie & Cooley

↳ reduce the complexity of FT from $O(N^2)$ to $O(N \log N)$

using an observation \rightarrow correlation between coefficients \rightarrow break a problem of size N into 2 problems of size $\frac{N}{2}$

↳ break problem into N problem of size 1

the glue it together using an algorithm called butterfly alg.



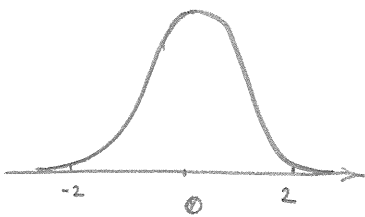
\rightarrow butterfly

↳ to fix it, we need to shift it, flip the sides and multiply other components by -1
or use `fftshift` and absolute value of it.

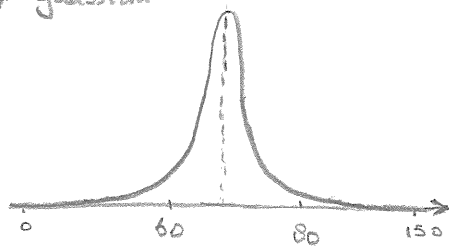
$$ut = \text{fft}(u)$$

$$\text{plot}(\text{abs}(\text{fftshift}(ut)))$$

✓ FT of a gaussian func is another gaussian



\Rightarrow



\rightarrow this axis is frequency
but we didn't order this

↳ matlab thinks we're working on 2π periodic domain
but we work on L period!

and we should rescale $\rightarrow \frac{2\pi}{L}$

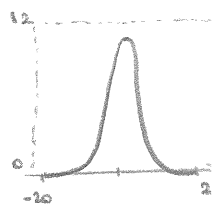
↳ $[0: \frac{n}{2}-1 \quad -\frac{n}{2}: -1]$ \rightarrow shifted integers (because fft is shifted too)

so

$$k = \frac{2\pi}{L} \times [0: \frac{n}{2}-1 \quad -\frac{n}{2}: -1]$$

$$\text{plot}(\text{fftshift}(k), \text{abs}(\text{fftshift}(ut)))$$

\Rightarrow Spectral Content of u



- Spectral Content

↳ how much energy is at cosine ωx ? ...

↳ tells us the strength of all of the freq. content

- always have it in mind that it is shifted

• Derivative Relations

$\widehat{f^{(n)}} = (ik)^n \widehat{f}$

↳ $\widehat{f^{(n)}}$ represents Fourier transform
 ↳ \widehat{f} represents n^{th} derivative

- to find 3rd derivative of f

1. Fourier transform it
2. multiply it by $(ik)^3$
3. inverse transform the result!

- if the boundary conditions are satisfied, this is one of the most accurate ways to calc derivatives

$u_j d s = \text{ifft}((i * k).^j .* u_t);$ → j^{th} derivative of function u using spectral approximation

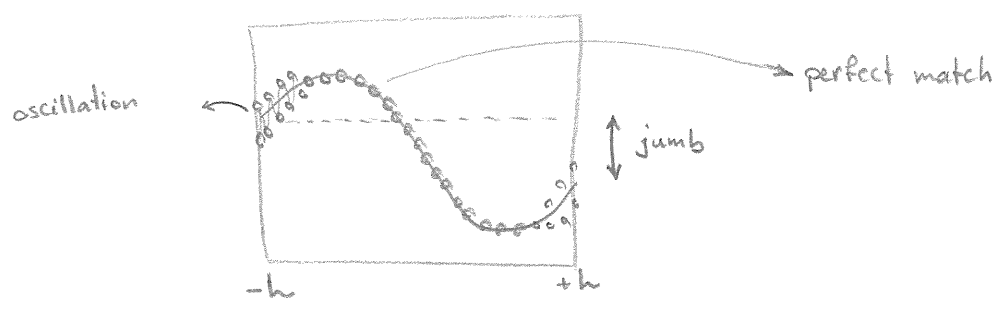
↳ it's accurate to the 10^{-6} magnitude

- if the domain is smaller, e.g. $h=20 \rightarrow h=4$

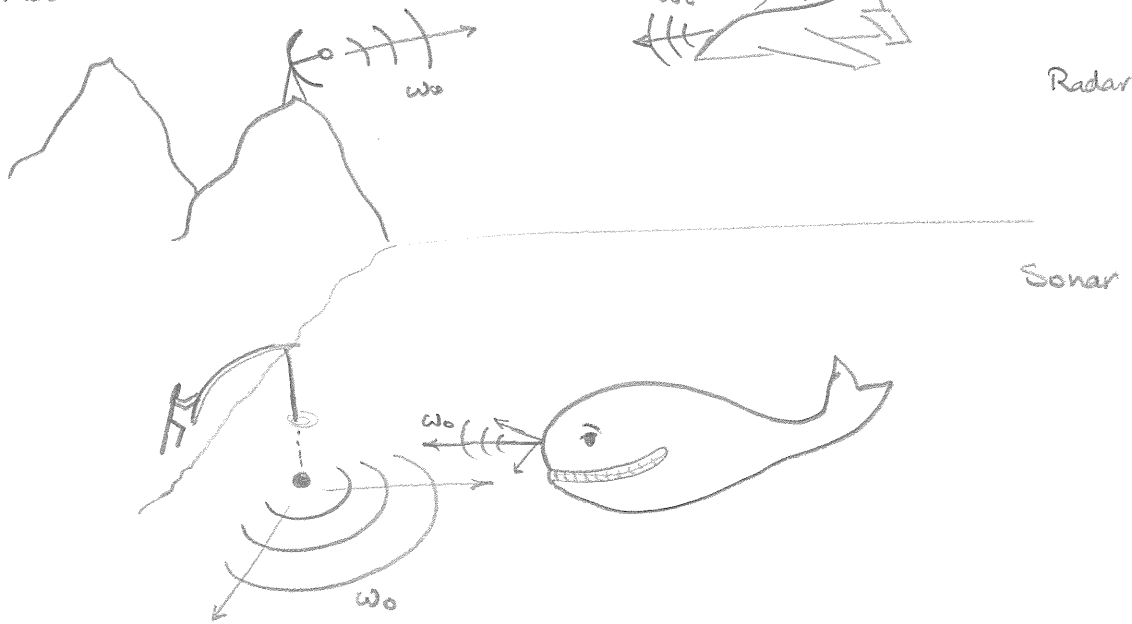
↳ violates the conditions of boundary

↳ it supposed to be 2π periodic function, but there will be a jump between the value of $u(-L)$ and $u(L)$ if h is small.

↳ result has a "Gibbs's Phenomenon" = oscillations near the jump



* Radar Problem

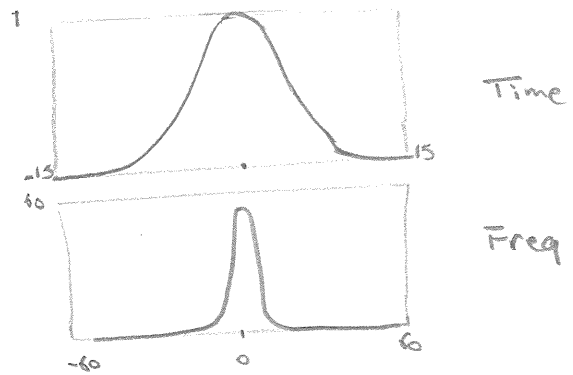


- Transmitter sends signal with ω_0
- Jets try to absorb that signal to be stealth
- Planes reflect a signal with freq. ω_0
- AVAX planes flood the area with ω_0 freq to cover the jets from being identified
- There are all lots of electromagnetic noise

• Heisenberg Uncertainty Principle

- there's a relationship between the width of the pulse in time and the width of the pulse in frequency.

$$\begin{aligned} &\rightarrow \text{width of time} \\ &\times \text{width of freq} \\ &= \\ &2\pi \end{aligned}$$



↳ Quantum Mechanics was built on this idea

↳ you can represent the state/function of a particle in terms of its probability density and a FT of that represent its momentum.

↳ you can never know the momentum and position at the same time because there's an uncertainty relationship between them

↳ width of the momentum and width of the location have to be related ($= 2\pi$)

□ narrow signal in time \rightarrow wide in freq: lots of freq components
broad in time \rightarrow becomes almost a single frequency

• white Noise

- white light = collection of all colors
- white noise = collection of all frequencies

* State-Space Model

- latent Markov process \rightarrow no time dependant X_t
- observation Y_t
- HMM $X_t = \Psi(X_{t-1}, V_t)$ $Y_t = \Phi(X_t, W_t)$
- Aim: infer X_t given observations Y_t

• Inference

$$X_1 \sim \mu \quad X_t | (X_{t-1} = x) \sim f(\cdot | x)$$

$$Y_t | (X_t = x) \sim g(\cdot | x)$$

$$P(y_{1:t}) = \int \dots \int P(x_{1:t}, y_{1:t}) dx_{1:t}$$

• MC Basics

$$\hat{P}(x_{1:t} | y_{1:t}) = \frac{1}{N} \sum_{i=1}^N \delta_{x_{1:t}}^{(i)}(x_{1:t})$$

- integration & marginalization is easy

- problem \rightarrow sampling from non-Gaussian non-linear models
 \searrow even if we can, it takes $O(t)$

- solution \rightarrow use MCMC for sampling approximately

- SMC Methods \rightarrow break problem of sampling to easier one

• Bayesian Recursion

$$P(x_t | y_{1:t-1}) = \int P(x_{t-1}, x_t | y_{1:t-1}) dx_{t-1}$$

$$= \int f(x_t | x_{t-1}) p(x_{t-1} | y_{1:t-1}) dx_{t-1}$$

مسائل سنگین - کار کثرت

• Bayesian Recursion on Path Space

* typos

$$P. 47 \rightarrow P(x_t | y_t, x_{t-1}) = \frac{f(x_t) g(y_t | x_t)}{P(y_t | x_{t-1})}$$

Submodularity 8

(Bach - Inria)

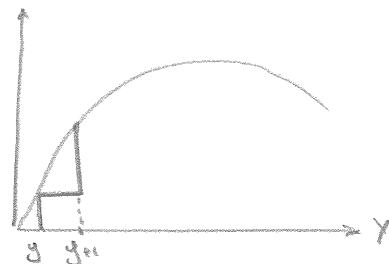
* Problem

• F is submodular iff $F: 2^V \rightarrow \mathbb{R}$

$$\forall A, B \subseteq V, F(A) + F(B) \geq F(A \cup B) + F(A \cap B)$$

$F(\emptyset) = 0 \rightarrow$ submodular

• other definition



• enough to show that

take $A \rightarrow$ Add one item to A

$$\forall A, k, j \quad F(A \cup \{k, j\}) + F(A) \leq F(A \cup \{j\}) + F(A \cup \{k\})$$

• Closedness Properties

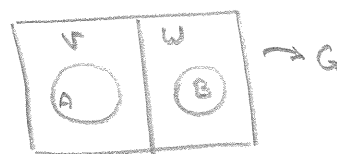
- A & B : submodular \rightarrow linear combination

Marginalization $A \mapsto F(A \cap B)$

Conditioning $A \mapsto F(A \cup B) - F(B)$

• Partial Minimization

$$F(A) = \min_B G(A \cup B) - \min_B G(B)$$



• Modular Function

$$S(A) = S^T I_A \quad \text{indicator vector of } A$$

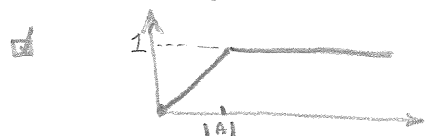
$$S(A) = \sum_{k \in A} S_k$$

$$S = I_A \Rightarrow S(A) = |A|$$

- $F: A \mapsto g(S(A))$ is submodular if g is a concave function

$$F(A) = g(|A|)$$

proof: $g(|A|+2) + g(|A|) \leq g(|A|+1) + g(|A|+1) \Rightarrow g(k+1) - g(k) \geq g(k+2) - g(k+1)$



* Matroids

- Essentially a column of matrix \rightarrow independent sets
- Rank function \rightarrow biggest sub matrix in matroid that has independence

* Technical Challenges

- Put object on screen correctly in its position
- How to draw surfaces
- How to light an object correctly
 - Highlights on it
 - Smooth shading
- How are objects lit in real world?
- Build Softwares and Hardwares

* 3D Graphics Pipeline

1. Modeling: Creating geometric models of object
 - Simple object Teapot
 - Complicated Mesh Cat Sculpture
 - Million Polygon Mesh Michelangelo's David
2. Animation: Motion of character
3. Rendering: Creating realistic images, given the geometry and animation
 - Simulate the way light propagate on scene → realistic & intricate shadows

* Rasterize vs Ray Trace

- Rasterization: Goes through all the geometric primitives and determines where in the screen they should go
- Ray Tracing: Goes to each point/pixel in the screen and determines which geometric primitive that corresponds to.



RECEIPT

Paid by : Meshgi Kourosch
(meshgi-k@sys.i.kyoto-u.ac.jp)
Kyoto University

Paid to : Machine Learning Summer School'12 in Kyoto

MLSS12 Secretariat Office,
Graduate School of Informatics, Kyoto University
Research Bldg. No.1/Project Lab, Room 303,
36-1 Yoshida-Honmachi, Sakyo-ku, Kyoto 606-8501 JAPAN
+81 (0)75-753-5911

The sum of: 40,000 JPY

Date: August 23, 2012

corresponding to:

Description	Quantity	Unit price	Amount
Registration for the MLSS'12 in Kyoto, August 27 - September 7 2012 - Student, 2 weeks	1	40,000 JPY	40,000 JPY

Total: 40,000 JPY (Tax: 0 JPY)

Terms and conditions

The registration includes: access to all the lectures publicized on the MLSS schedule, light snacks during coffee breaks between lectures, banquet and MLSS party event described in the "Social Events" section of the MLSS website.

機 械 学 習
サマースクール
京 都 2012

	MON. 27	TUE. 28	WED. 29	THU. 30	FRI. 31
8:30 - 10:10	Opening	Domingos <i>Stat. Rel. Learn.</i>	Vandenberghes <i>Convex Optim.</i>	Vandenberghes <i>Convex Optim.</i>	Lin <i>ML Software</i>
10:30 - 12:10	Rakhliln <i>Stat. Learning</i>	Rakhliln <i>Stat. Learning</i>	Vandenberghes <i>Convex Optim.</i>	Müller <i>Brain C. I.</i>	Lin <i>ML Software</i>
	Lunch Break				
13:50 - 15:30	Rakhliln <i>Stat. Learning</i>	Tsuda <i>Graph Mining</i>	Tsuda <i>Graph Mining</i>	Müller <i>Brain C. I.</i>	Schapire <i>Boosting</i>
15:50 - 17:30	Domingos <i>Stat. Rel. Learn.</i>	Tsuda <i>Graph Mining</i>	Müller <i>Brain C. I.</i>	Schapire <i>Boosting</i>	Schapire <i>Boosting</i>
17:50 - 19:30	Domingos <i>Stat. Rel. Learn.</i>	Poster I <i>Intern. Hall I&II</i>	Doya <i>Brain & R. L.</i>	Poster II <i>Intern. Hall I&II</i>	Okada <i>ML & Med. Im.</i>
	MON. 3	TUE. 4	WED. 5	THU. 6	FRI. 7
8:30 - 10:10	Wainwright <i>Graphical Models</i>	Blei <i>Topic Models</i>	Blei <i>Topic Models</i>	Vempala <i>High-dim. Sampling</i>	Fukumizu <i>Kernel Methods</i>
10:30 - 12:10	Wainwright <i>Graphical Models</i>	Blei <i>Topic Models</i>	Vempala <i>High-dim. Sampling</i>	Fukumizu <i>Kernel Methods</i>	Fukumizu <i>Kernel Methods</i>
	Lunch Break				
13:50 - 15:30	Doucet <i>Seq. Monte Carlo</i>	Doucet <i>Seq. Monte Carlo</i>	Vempala <i>High-dim. Sampling</i>	Bach <i>Submodularity</i>	Bach <i>Submodularity</i>
15:50 - 17:30	Doucet <i>Seq. Monte Carlo</i>	Wainwright <i>Graphical Models</i>	Takemura <i>Holonomic Grad.</i>	Bach <i>Submodularity</i>	Sugiyama <i>Density Ratio</i>
17:50 - 19:30	Poster III <i>Intern. Hall I&II</i>	Amari <i>ML & Info. Geom.</i>	Banquet	Iwata <i>Submodular Optim.</i>	