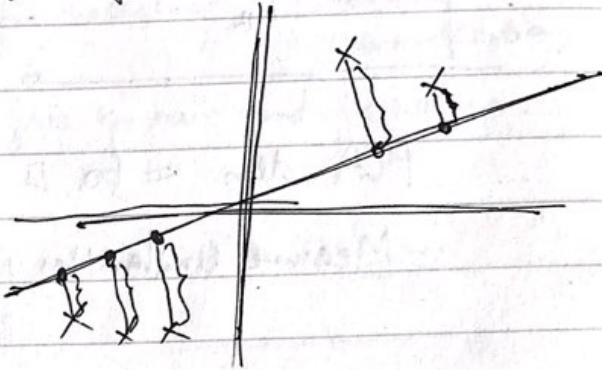


Plane $x^{(i)} \in \mathbb{R}^n$. View representation of data
in $\{u_1, \dots, u_k\}$ basis.

$$y^{(i)} = (u_1^T x^{(i)}, u_2^T x^{(i)}, \dots, u_k^T x^{(i)}). y^{(i)} \in \mathbb{R}^k$$

$$\min \sum \{\text{distance}\}^2$$

More interpretations:
See Problem Set.



Applications of PCA

- Visualization

Krishna-Shenoy's Lab.

- Compression

- Learning

less prone to overfitting

over used in industry

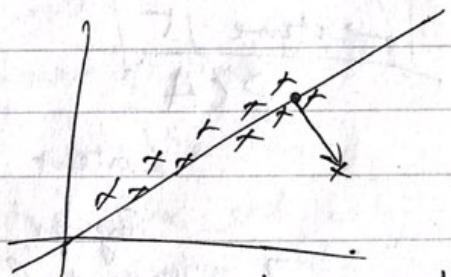
- Anomaly detection

- Matching / distance calculations

(100x100 image)

(face detection)

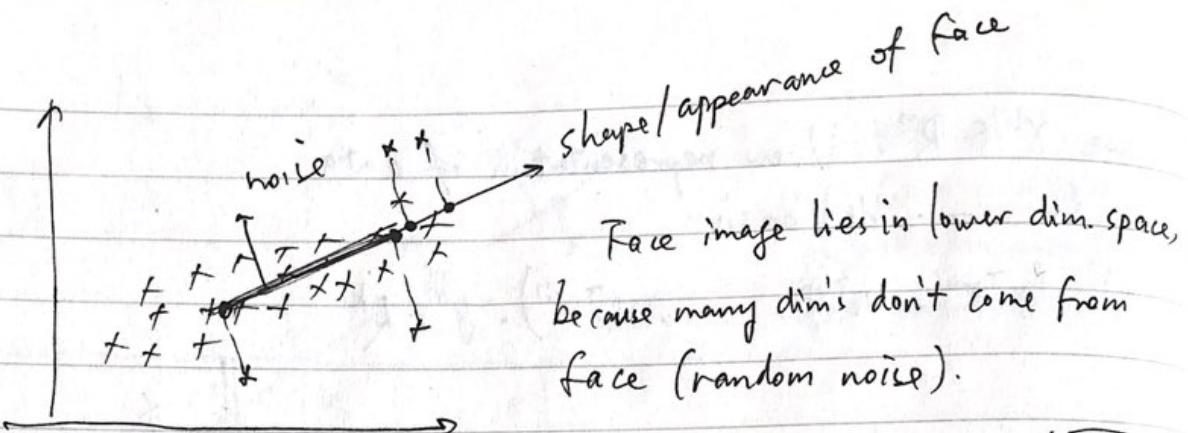
$$x^{(i)} \in \mathbb{R}^{10000}$$



Not the best anomaly
detection algorithm, but
sometimes it is enough.

○	○	○
○	○	○
○	○	○
○	○	○
○	○	○

"Sandy Alex Penland's Lab @ MIT"



PCA dim. = f_0 is typical.

Measure similarities in subspace.

$$x^{(i)} \in \mathbb{R}^{10000} \rightarrow \boxed{\text{♀}}$$

$$u^{(i)} \in \mathbb{R}^{10000} \rightarrow \boxed{\text{♀}}_{\text{eigenface}}$$

$x^{(i)}$ is linear combinations of eigenfaces

* Often recommend PCA. [Andrew Ng].

But over-used in industry (more often used):

* Before using PCA, think whether need compression exactly.

Lecture 15.

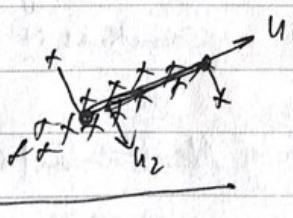
PCA

- Latent Semantic Indexing (LSI)

- Singular Value Decomposition (SVD) implementation

Independent Component Analysis (ICA)

PCA



Algorithm:

1) Normalize to zero mean, unit variance

$$3) \Sigma = \frac{1}{m} \sum_{i=1}^m x^{(i)} x^{(i)\top}$$

3) Find top k eigenvectors of Σ .

$$x^{(i)} \in \mathbb{R}^{10000} \quad 100 \times 10000$$

$$\Sigma \in \mathbb{R}^{10000 \times 10000}$$

?

PCA: Text data

$$\begin{bmatrix} 1 \\ 0 \\ \vdots \\ 1 \\ 0 \\ \vdots \end{bmatrix} \begin{array}{l} \text{ardvark} \\ | \\ \text{learn} \\ | \\ \text{study} \\ | \end{array}$$

LSI

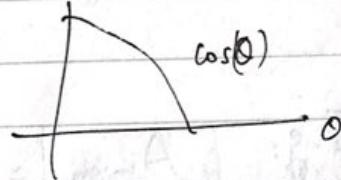
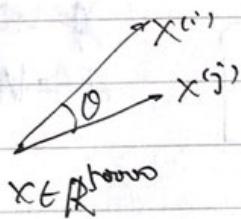
(Latent Semantic Indexing)

(Usually skip preprocessing: variance normalization)

if do that, weights of rare word will scale up.

$x^{(i)}, x^{(j)}$ - measure similarity.

$$\text{Sim}(x^{(i)}, x^{(j)}) = \cos \theta = \frac{x^{(i)T} x^{(j)}}{\|x^{(i)}\| \cdot \|x^{(j)}\|}$$



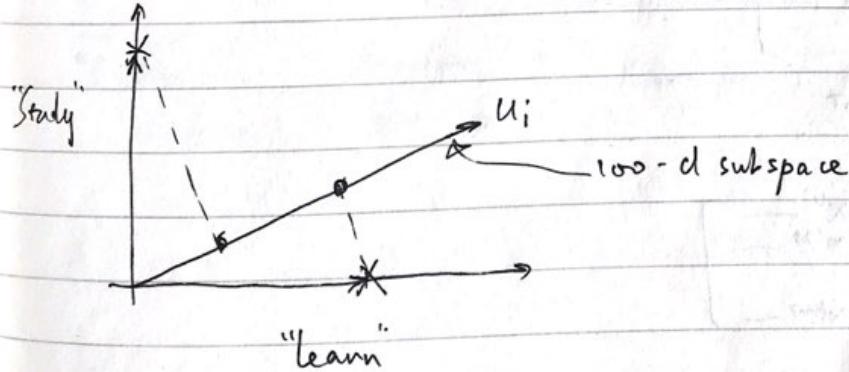
Eg: $x^{(i)} \rightarrow \text{"study"}$

$x^{(j)} \rightarrow \text{"learn"}$

Numerator of $\text{Sim}(x^{(i)}, x^{(j)})$

$$x^{(i)T} x^{(j)} = \sum_k x_k^{(i)} x_k^{(j)}$$

$$= \sum_k \mathbb{1}\{\text{Doc } i \text{ and } j \text{ both contain word } k\}$$



$\text{Sim}(x^{(i)}, x^{(j)}) > 0$ after PCA.

Implementation PCA

$$X^{(i)} \in \mathbb{R}^{50000}$$

$$\Sigma \in \mathbb{R}^{50000 \times 50000}$$

$$A \in \mathbb{R}^{n \times n}$$

$$A = U D V^T \quad (\text{SVD})$$

$$D = \text{diagonal} = \begin{bmatrix} \sigma_1 & & \\ & \sigma_2 & \\ & & \ddots \\ & & & \sigma_n \end{bmatrix} \quad \sigma_i = \text{singular values of } A.$$

$$\text{E.g.: } [A_{m \times n}] = [U_{m \times n}] [D_{n \times n}] [V^T_{n \times n}]$$

$$\cancel{A} = UDV^T$$

$$ATA = UDV^T U^T VD^T V^T$$

* U's columns: Eigenvectors of AAT^T ✓

* V's columns: Eigenvectors of ATA^T ✓

Compute using "svd" in Matlab/Octave (Very stable and efficient)

$$[A_{m \times n}] = [U_{m \times n}] [D_{n \times n}] [V^T_{n \times n}]$$

$$\Sigma = \sum_{i=1}^n X^{(i)} X^{(i)T}.$$

Design Matrix $X = \begin{bmatrix} \underline{\underline{x^{(1)}}} \\ \underline{\underline{x^{(2)}}} \\ \vdots \\ \underline{\underline{x^{(m)}}} \end{bmatrix}$

$$\Sigma = X^T X = \begin{bmatrix} | & | & | \\ X^{(1)} & X^{(2)} & \cdots & X^{(m)} \\ | & | & | \end{bmatrix} \begin{bmatrix} \underline{\underline{x^{(1)}}} \\ \underline{\underline{x^{(2)}}} \\ \vdots \\ \underline{\underline{x^{(m)}}} \end{bmatrix}$$

To get top k eigenvectors

$$X = UDV^T$$

V? U?

Top k columns of V are top k eigenvectors of $X^T X = \Sigma$.

$$X \in \mathbb{R}^{m \times 50000}$$

$$\Sigma \neq XX^T$$

A note of caution:

$$[X_{mn}] = [U_{mn}] \begin{bmatrix} D & \\ & 0 \end{bmatrix} [V^T]$$

$$[U_{mn}] \begin{bmatrix} D & \\ & 0 \end{bmatrix} [V^T]$$

dim of matrix

"Conventions" may
be different between diff.
SVD commands.

	Model $P(x)$	Not probabilistic
"Subspace"	Factor Analysis	PCA
"Clumps"/ "Groups"	Mixture of Gaussians	K-means

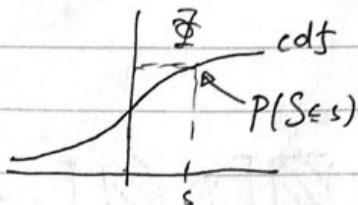
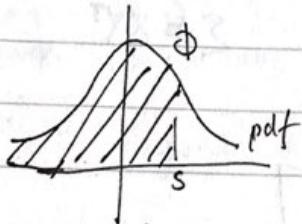
ICA

cdf (cumulative distribution function)

random variable S , density $P_S(s)$

cdf: $F(s) = P(S \leq s)$

E.g:



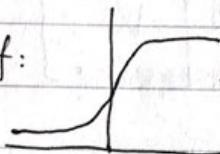
$$F(s) = \int_{-\infty}^s P_S(t) dt$$

\rightarrow Specify $P_S(s)$

or Specify $F(s)$

$$P_S(s) = F'(s)$$

cdf:



Sigmoid

Any function:

- continuous

- $f(-\infty) \approx 0$

- $f(+\infty) = 1$

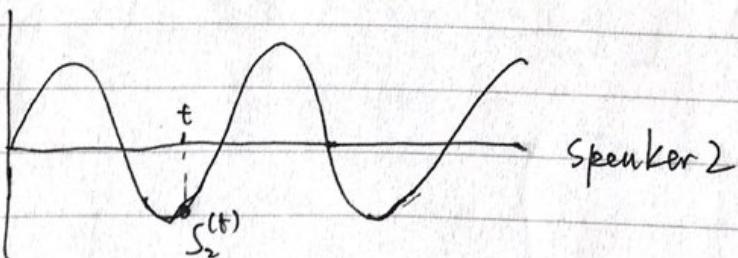
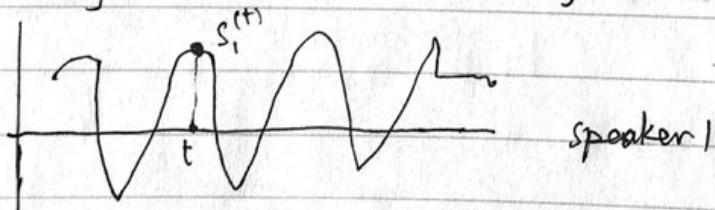
can do.

Cocktail Party

Original sources: n speakers

$$S \in \mathbb{R}^n$$

$S_j^{(t)}$ = signal from speaker j at time t .



We observe : n microphones

$$x^{(i)} = As^{(i)} \quad x^{(i)} \in \mathbb{R}^n$$

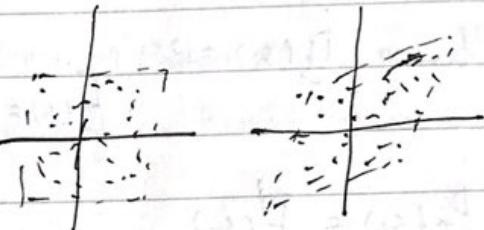
$$x_j^{(i)} = \sum_k A_{jk} s_k^{(i)}$$

$$W = \begin{bmatrix} w_1^\top \\ w_2^\top \\ \vdots \\ w_n^\top \end{bmatrix}$$

Goal: Find $W = A'$ so that
 $s^{(i)} = Wx^{(i)}$

If $s_j^{(i)} \sim \text{Uniform } [-1, 1]$

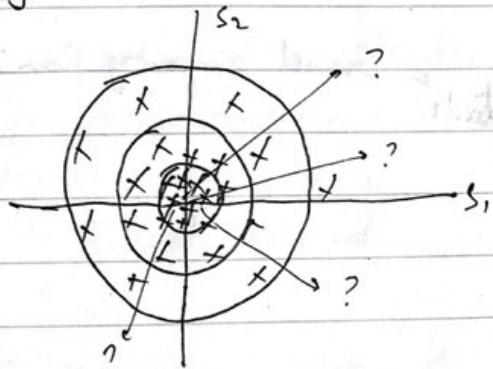
Ambiguities { Permutation of speaker
 | Sign of speaker signal }



$s_j^{(i)}$ - non-Gaussian

$$s^{(i)} \sim N(0, I)$$

If sources is Gaussian, it's
actually impossible to do ICA !!



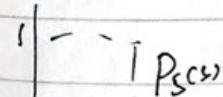
Let $s \in \mathbb{R}^n$.

Density of s : $P_s(s)$

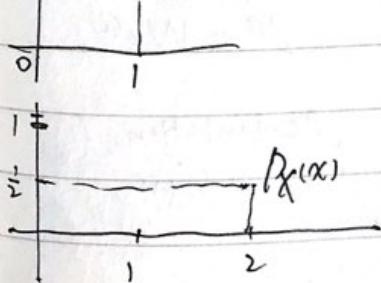
$$x = As = W^{-1}s, \quad s = Wx$$

$$P_X(x) = P_s(Wx) \cdot |W|$$

E.g.: $P_s(s) = \mathbb{I}\{0 \leq s \leq 1\}$ ($s \sim \text{Uniform } [0, 1]$).



Let $x = 2s$, $A = 2$, $W = \frac{1}{2}$



$$P_X(x) = \mathbb{I}\{0 \leq x \leq 2\} \cdot \frac{1}{2}$$

ICA model

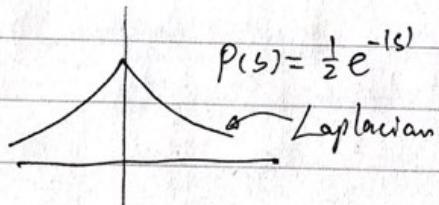
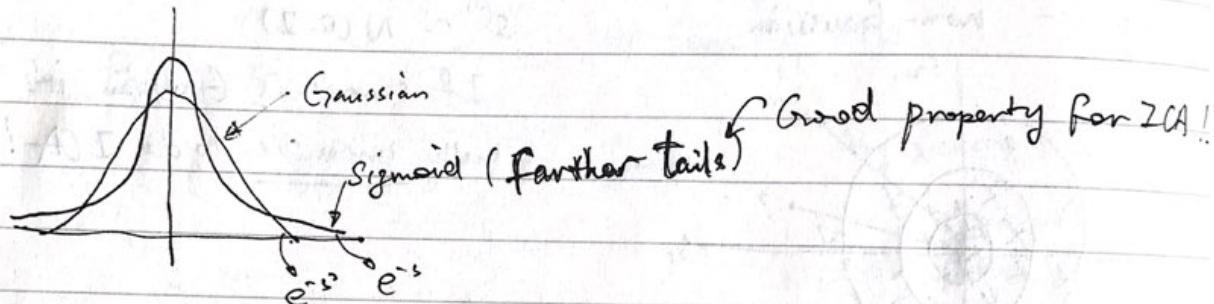
$$P(s) = \prod_{i=1}^n P_S(s_i) \quad (s_i's \text{ independent})$$

$$P(x) = \left[\prod_{i=1}^n P_S(w_i^T x) \right] / |W| \quad W = A^{-1} = \begin{bmatrix} w_1^T \\ w_2^T \\ \vdots \\ w_n^T \end{bmatrix}$$

Choose $P_S(s_i) = ?$ $s_i = w_i^T x.$

$$F(s) = \frac{1}{1+e^{-s}} \quad \text{Convenience, familiar -}$$

$$P_S(s_i) = F'(s_i)$$



Given $\{x^{(1)}, \dots, x^{(m)}\}.$

$$\ell(w) = \sum_i \log \left[\left(\prod_j P_S(w_j^T x^{(i)}) \right) / |W| \right] \quad P_S(s) = F'(s)$$

Stochastic gradient ascent:

$$w = w + \alpha \nabla_w \ell(w)$$

$$\nabla_w \ell(w) = \begin{bmatrix} 1 - 2g(w_1^T x^{(1)}) \\ \vdots \\ 1 - 2g(w_n^T x^{(1)}) \end{bmatrix} X^{(1)T} + (w^T)^{-1} \quad \begin{cases} \text{Output} \\ s^{(1)} = w x^{(1)} \end{cases}$$

Permutation
+/- sign

Applications of ICA

① Katie (T.A. of Ng) Photo: electrons on his bread.

EEG cap: pre-processing for research on brain.

② Aapo

Natural Images (~~pieces~~ Image Piece \rightarrow Pattern Piece).

hypothesis on early visual processing in human brain

Lecture 16

Reinforcement Learning \leftarrow "Between supervised Learning & unsupervised Learning"

- Markov Decision Process (MDP)
 - Value Function
 - Value Iteration
 - Policy Iteration
- "Much less supervision"
- "Just like training a dog"

Not one-time decision, but sequential decision process.

Credit Assignment Problem: which step cause to bad result?

MDP: Markov Decision Process

MDP $(S, A, \{P_{sa}\}, \mathcal{V}, R)$ "five-tuple"

① S - set of state (position & orientation)

② A - set of actions (position of control stake)

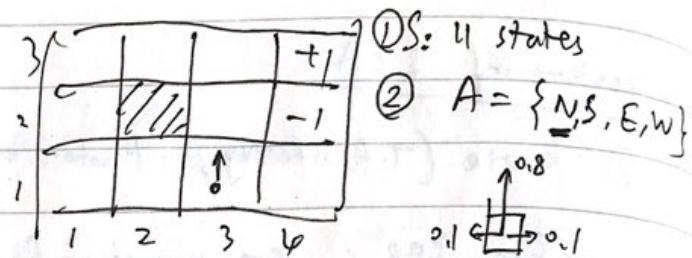
③ P_{sa} - state transition distributions

$$\sum_{s'} P_{sa}(s') = 1 . \quad P_{sa} \geq 0 .$$

④ \mathcal{V} - discount factor $0 \leq \mathcal{V} \leq 1$

⑤ R - reward function $R: S \mapsto \mathbb{R}$

eg: from AI: A Modern Approach
Russell & Norvig



$$③ P_{(3,1)N}((3,2)) = 0.8$$

At state s_0

$$P_{(3,1)N}((4,1)) = 0.1$$

Choose a_0

$$P_{(3,1)N}((2,1)) = 0.1$$

Get to $s_1 \sim p_{s_0, a_0}$

$$P_{(3,1)N}((3,1)) = 0$$

Choose a_1

$$\vdots$$

$$\vdots \text{ etc.}$$

Get to $s_2 \sim p_{s_1, a_1}$

$$④ R((4,3)) = +1$$

$$R((4,2)) = -1$$

$$R(s) = -0.02 \text{ for all other states: "Battery charges and time consume"}$$

⑤ Total Payoff:

$$R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \dots$$

$$0 \leq \gamma < 1 \text{ (often } 0.99)$$

In economic app.

represent interest of bank etc..

Win / lost immediately greater than future.

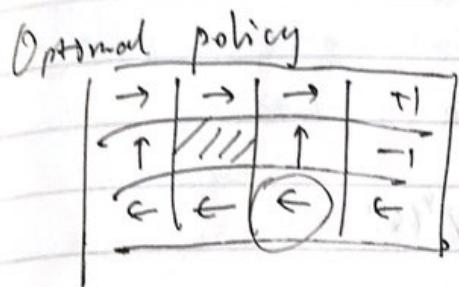
* Choose actions over time (a_0, a_1, \dots) to maximize

$$E [R(s_0) + \gamma R(s_1) + \gamma^2 R(s_2) + \dots]$$

* Policy $\pi: S \rightarrow A$

Policy: Simple, only depend on current time

Strategy: Multiple policies, include previous states



MDP good at:

Subtle tradeoff between
different impact factors

(Conservative policy \leftrightarrow Active policy)
"low risk" "smaller distance"

Define V^π , V^* , π^*

For any π , define value function $V^\pi : S \mapsto \mathbb{R}$

st. $V^\pi(s)$ is expected total payoff starting in state s ,
and executed π .

$$V^\pi(s) = E [R(s_0) + \gamma R(s_1) + \dots | \pi, s=s_0]$$

e.g.:

$$V^\pi(s) = \sum_{s'} P_{\pi(s)}(s') V^\pi(s')$$

$V^\pi(s) = R(s) + \gamma \sum_{s'} P_{\pi(s)}(s') V^\pi(s')$

Bellman's Equations

$$V^\pi(s) = ?$$

$$\underline{V^\pi((3,1))} = R((3,1)) + \gamma [0.8 \underline{V^\pi((3,2))} + 0.1 \underline{V^\pi((4,1))} + 0.1 \underline{V^\pi((2,1))}]$$

$$\pi((3,1)) = N \begin{array}{|c|c|c|} \hline & \checkmark & \\ \hline \checkmark & & \checkmark \\ \hline \end{array}$$

Optimal value function:

$$V^*(s) = \max_{\pi} V^\pi(s)$$

Bellman's eqns:

$$\left\{ \begin{array}{l} V^*(s) = R(s) + \max_a \sum_{s'} P_{sa}(s') V^*(s') \end{array} \right.$$

$$\left\{ \begin{array}{l} \pi^*(s) = \arg \max_a \sum_{s'} P_{sa}(s') V^*(s') \end{array} \right.$$

Value iteration ("individual", "liberalism")

Initialize $V(s) = 0$ for all s .

For every s , update

$$V(s) := R(s) + \max_a \sum_{s'} P_{sa}(s') V(s')$$

This will make $V(s) \rightarrow V^*(s)$

Synchronous update

$$V := B(V)$$

Bellman backup operator

Asynchronous update

(a little bit faster)

.86	.90	.93	+1
.82	/ /	.69	-1
.78	.75	.71	.49

V^*

$$\pi((3,1)) = w : \sum_{s'} P_{3a}(s') V^*(s')$$

$$= .8 \times .75 + .1 \times .69 + .1 \times .71 = .74$$

Stay when hit wall

→	→	→	+1
↑	/ /	↑	-1
←	←	←	←

π^*

$$\pi((3,1)) = n : \sum_{s'} P_{3a}(s') V^*(s')$$

$$= .8 \times .69 + .1 \times .75 + .1 \times .49 = .676$$

Policy Iteration ("Centralised")

Initialize π randomly.

Repeat:

Let $V := V^\pi$ (solve Bellman's eq^{ns}) (e.g.: 10 million) -

Let $\pi(s) := \arg \max_a \sum_{s'} P(s'|s) V(s')$ Better choose value iter.

$V \rightarrow V^*$ $\pi \rightarrow \pi^*$ (At last)

Also Convex optimization

Proof not hard

What if don't know P_{sa} ?

estimate P_{sa} from data

$S, A, \{P_{sa}\}, V, R$

$\checkmark \quad \checkmark \quad ? \quad \checkmark \quad \times$

$$P_{sa}(s') = \frac{\text{\# times took action "a" in } S, \text{ got to } s'}{\text{\# times took action "a" in } S}$$

(or $\frac{1}{|S|}$ if "0") uniform distribution when no experiments

More mix & match algo's in R.L.

Putting it together

repeat:

Take actions using π to get experience in MDP

Update estimates of P_{sa}

Solve Bellman's eq^{ns} (V^*) using V.I. to get V^*

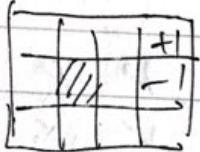
Update $\pi(s) = \arg \max_a \sum_{s'} P_{sa}(s') V^*(s')$

}

Lecture 17

- Continuous State MDPs
- Discretization
- Models / Simulators
- Fitted Value Iteration
- Q function
- Approximate policy iteration

MDP $\{S, A, \{P_{sa}\}, V, R\}$



S : 11 states

$A = \{N, S, E, W\}$

$P_{sa}:$

0.8	
0.1	0.1

$R = \begin{cases} \pm 1 & \text{wt. absorbing states} \\ -0.02 & \text{elsewhere} \end{cases}$

$\gamma = 0.99$

$\pi : S \rightarrow A$

$$V^\pi(s) = E [R(s_0) + \gamma R(s_1) + \dots | \pi, s=s_0]$$

Find $V^*(s) = \max_\pi V^\pi(s)$

$$\pi^*(s) = \arg \max_\pi \sum_{s'} P_{sa}(s) V^*(s')$$

$$V^*(s) = R(s) + \gamma \max_a \sum_{s'} P_{sa}(s') V^*(s')$$

$$V.I.: V(s) := R(s) + \gamma \max_a \sum_{s'} P_{sa}(s') V(s')$$

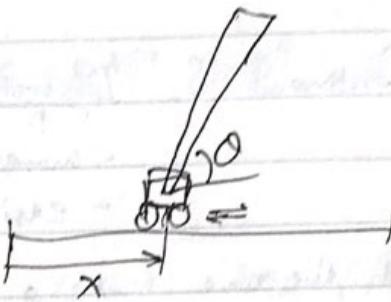
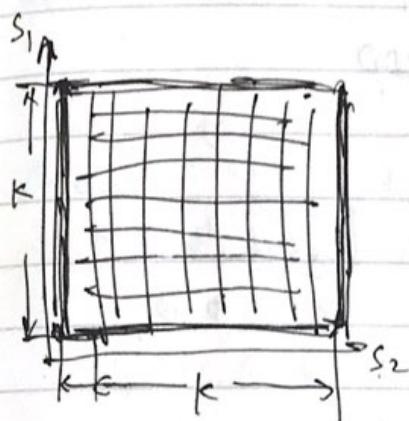
P.I.: Repeat { Solve for V^π . Let $V := V^\pi$

$$\text{Update } \pi(s) = \arg \max_a \sum_{s'} P_{sa}(s') V(s')$$

E.g.: Car: $x, y, \theta, \dot{x}, \dot{y}, \dot{\theta}$

Helicopter: $x, y, z, \phi, \theta, \psi, \dot{x}, \dot{y}, \dot{z}, \dot{\phi}, \dot{\theta}, \dot{\psi}$
(roll, pitch, yaw)

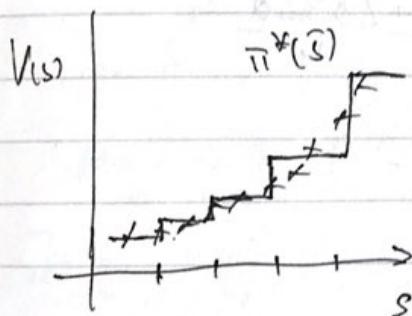
Inverted pendulum: $x, \theta, \dot{x}, \dot{\theta}$



Discretize into discrete states \bar{s}

Solve for $V^*(\bar{s})$, $\pi^*(\bar{s})$

execute $\pi^*(\bar{s})$



Not really good.

Not smooth between buckets

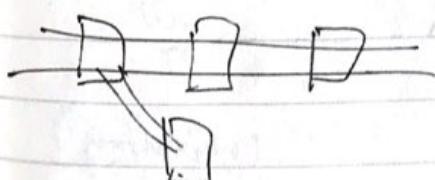
Not generalize between buckets

Curse of dimensionality.

If $s \in \mathbb{R}^n$ $n \leq 3$

Get k^n states. "folk trick"

Factory Automations



Scheduling orders

Board Games (e.g. play chess)

N machines $\times k$ states/machine

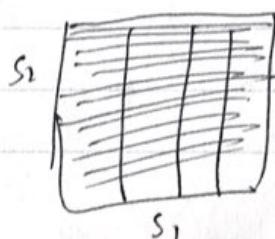
$= k^N$ states.

N pieces $\times k$ position/pieces

$= k^N$ states (positions).

discretization scale poorly for
high-dimensional spaces

How: Approximate V^* !



sensitive difference
for diff. dim.

manually HARD

Continuous S , Discrete A

more rational

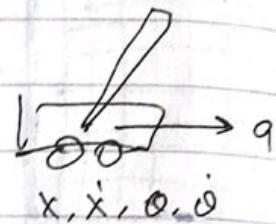
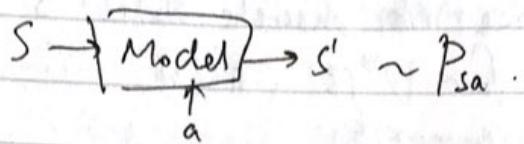
easier

car: 6 states vs. 2 actions

helicopter: 12 states vs. 4 actions

Inv.P.: 4 states vs. 1 action

① Assume have a model/simulator of MDP.



E.g.: physics simulator

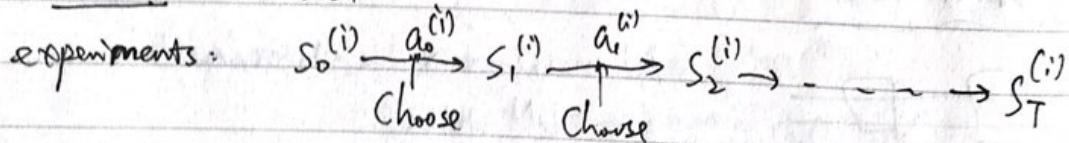
$$S = \begin{pmatrix} x \\ \dot{x} \\ \theta \\ \dot{\theta} \end{pmatrix}$$

$$\dot{S} = \begin{pmatrix} \dot{x} \\ \ddot{x} - L\beta \cos\theta/m \\ \dot{\theta} \\ \ddot{\theta} \end{pmatrix}, \text{ where } \begin{cases} \alpha = \frac{a + L\dot{\theta}^2 \sin\theta}{M} \\ \beta = \dots \\ \text{etc.} \end{cases}$$

$$S_{t+1} = S_t + f(\Delta t) \begin{pmatrix} \dot{x} \\ \ddot{x} - L\beta \cos\theta/m \\ \dot{\theta} \\ \ddot{\theta} \end{pmatrix}$$

$$\Delta t = 0.1 \text{ second}$$

② Learn a model



(Repeat i times).

"trajectory"

Use learning algo. to estimate S_{T+1} as function of S_t, a_t .

[Inv.P.: $S_t \in \mathbb{R}^4, a_t \in \mathbb{R}^4$. run L.R. 4 times]

E.g.: $S_{t+1} = A S_t + B a_t$ (linear), $A \in \mathbb{R}^{k \times k}$, $B \in \mathbb{R}^{k \times 4}$

$$\arg \min_{A, B} \sum_{i=1}^m \sum_{t=0}^{T-1} \|S_t^{(i)} - (AS_t^{(i)} + BA_t^{(i)})\|^2$$

L.R ... LWR, kernel L.P.

Model: $S_{t+1} = AS_t + Ba_t$ (deterministic)

or $S_{t+1} = AS_t + Ba_t + \varepsilon_t$, $\varepsilon_t \sim N(0, \Sigma)$ (stochastic)

ε_t sample from Gaussian distribution

$$S_t \xrightarrow[\text{at } a_t]{\text{Model}} S_{t+1}$$

Inv. P: LWR, Nonlinear Algo.

Approximate V^* !

Just like LR:

Choose feature $\phi(s)$ for state s

Could $\phi(s) = s$ parameters

Approximate: $V(s) = \theta^\top \phi(s)$

" $h_\theta(x) = \theta^\top \phi(x)$ "

E.g. Inv. P:

$$\phi(s) = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ \dot{x}_1 \\ \ddot{x}_2 \\ \ddot{\theta} \cdot x \\ \dot{\theta}^2 \\ \vdots \end{bmatrix}$$

V . I :

$$V(s) := R(s) + \gamma \max_a \sum_{s'} P_{sa}(s') V(s')$$

Replace by integral

$$= R(s) + \gamma \max_a E_{s' \sim P_{sa}} [V(s')]$$

Fitted Value Iterations

Sample $\{s^{(1)}, s^{(2)}, \dots, s^{(m)}\} \subseteq S$ randomly.

Initialize $\theta := 0$ (Just like L.R.)

Repeat {

For $i = 1 \dots m$ {

For each action $a \in A$ {

Sample $s'_1, \dots, s'_{k'} \sim P_{s^{(i)} a}$ (using model)

$$\text{Let } g(a) = \frac{1}{k'} \sum_{j=1}^{k'} [R(s^{(i)}) + \gamma V(s'_j)]$$

// estimate for $R(s^{(i)}) + \gamma E_{s' \sim P_{s^{(i)} a}}[V(s')]$

$$\text{Set } y^{(i)} = \max_a g(a)$$

// estimate for $R(s^{(i)}) + \gamma \max_a E_{s' \sim P_{s^{(i)} a}}[V(s')]$

}

// Want $V(s^{(i)}) \approx y^{(i)}$

$$\theta := \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n (\theta^T \phi(s^{(i)}) - y^{(i)})^2$$

}

helicopter control
collide - 100Hz
 ~ 100 samples

$$MDP = \{S, A, \{P_{sa}\}, R, R\}$$

Inr.-P. $\overset{\uparrow}{R^*}$ $\overset{\uparrow}{R}$ Model $\overset{\uparrow}{\pi_{\theta}}$

$$R(s) = \begin{cases} -1 & \text{fallen}/\theta \geq 30^\circ \\ 0 & \text{otherwise} \end{cases}$$

All above,

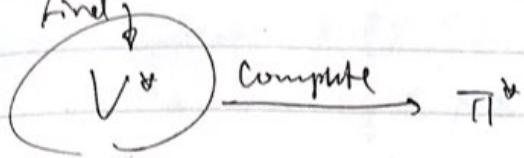
Assume: we have a stochastic simulator

Runs quite well
for 6-20 dim.
if choose features
properly.

If we have a deterministic simulator:

Set $k=1$. don't need sample (simplify F.V.I.)

Find



for L.R.: easy to choose feature.

For approx. value func.: harder/tricker.

discrete state s :

$$\pi^*(s) = \arg \max_a E_{s' \sim P_{sa}} [V^*(s')] = \arg \max_a \sum_s P_{sa}(s') V^*(s')$$

continuous state:

$$\text{At specific state } s = (x, \dot{x}, \ddot{x}, \ddot{\dot{x}})$$

$$a = \arg \max_a E_{s' \sim P_{sa}} [V^*(s')] \quad \text{Replace } s \rightarrow s'$$

If deterministic simulator:

$$\begin{array}{l|l} s_{t+1} = f(s_t, a_t) & \arg \max_a V^*(f(s, a)) \\ s' = f(s, a) & \end{array}$$

If stochastic simulator:

$$s_{t+1} = f(s_t, a_t) + \epsilon_t \quad \epsilon_t \text{ Gaussian noise}$$

$$\text{E.g.: } s_{t+1} = A s_t + B a_t + \epsilon_t$$

$$\begin{aligned} E_{s' \sim P_{sa}} [V^*(s')] &\approx V^*(E[s']) \\ &= V^*(f(s, a)) \\ &\quad \underbrace{\phi(s')}_{\phi \text{ from F.V.I.}} \end{aligned}$$

$$s' = f(s, a)$$

Choose action:

$$\arg \max_a V^*(f(s, a)) \quad \text{ignoring the noise (compute quickly)}$$

effective in practice (mean 0 Gaussian noise)
in model is proper

Lecture 18

State-action rewards

Finite horizon MDPs

Linear Dynamical Systems

- Models

- Linear Quadratic Regulation (LQR) control

- Riccati Equations

MDP: $(S, A, \{P_{sa}\}, \gamma, R)$ $0 \leq \gamma < 1$
 $R: S \mapsto \mathbb{R}$.

$$V(s) := R(s) + \max_a \gamma \sum_{s'} P_{sa}(s') V(s')$$

$$\pi^*(s) = \arg \max_a \gamma \sum_{s'} P_{sa}(s') V^*(s')$$

State-action Rewards

$R: S \times A \mapsto \mathbb{R}$.

~~Payoff~~

$s_0, a_0 \rightarrow s_1, a_1 \rightarrow s_2, a_2 \rightarrow \dots$

Payoff: $R(s_0, a_0) + \gamma R(s_1, a_1) + \gamma^2 R(s_2, a_2) + \dots$ / action "stay"

$$V^*(s) = \max_a R(s, a) + \gamma \sum_{s'} P_{sa}(s') V^*(s')$$

V.I: $V(s) := \max_a R(s, a) + \gamma \sum_{s'} P_{sa}(s') V^*(s')$

$$\pi^*(s) = \arg \max_a R(s, a) + \gamma \sum_{s'} P_{sa}(s') V^*(s')$$

Finite horizon MDPs.

$$(S, A, \{P_{sa}^{(t)}\}, T, R)$$

T: horizon time.

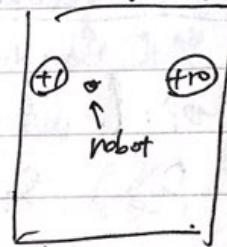
$$\text{Payoff: } R^{(0)}(s_0, a_0) + R^{(1)}(s_1, a_1) + \dots + R^{(T)}(s_T, a_T)$$

Optimal policy: non-stationary (depend on time)

Left OR right?

$$s_{t+1} \sim P_{s,a}^{(t)}$$

Spacecraft: fuel (time)
traffic: weather (time)
factory: workers (time).



depend on the ~~time~~ remained

$$V_t^*(s) = E[R^{(t)}(s_t, a_t) + \dots + R^{(T)}(s_T, a_T) | \pi^*, s_t = s]$$

$$V_T^*(s) = \max_a R^{(T)}(s, a)$$

$$V_t^*(s) = \max_a R^{(t)}(s, a) + \sum_{s'} P_{sa}^{(t)} V_{t+1}^*(s')$$

$$Tl_t^*(s) = \arg \max_a R^{(t)}(s, a) + \sum_{s'} P_{sa}^{(t)} V_{t+1}^*(s')$$

Dynamic Programming
(Recursion)

$$\left\{ \begin{array}{c} V_T^*, V_{T-1}^*, V_{T-2}^*, \dots, V_0^* \\ Tl_T^*, \dots, Tl_1^* \end{array} \right.$$

directly compute.
not iter. to converge.

$$V^{(T)} \approx 0.00001 \rightarrow T.$$

Commonly use discounting OR finite horizon.

① Do discounting can ensure that the value function is finite.

② Use finite horizon is sufficient to ensure the value function is finite.

Both use state-action rewards and finite horizon MDPs

(Strong assumption)

But it is reasonable for many systems.

Then we have a very elegant algo. to solve even every

large MDPs.

Linear Quadratic Regulation (LQR)

Continuous state space, maybe continuous action.

MDP: $\{S, A, \{P_{sa}\}, T, R\}$ $S \in \mathbb{R}^n$ $A \in \mathbb{R}^{d \times d}$

$$P_{ca}: S_{t+1} = A_s S_t + B_t a_t + w_t \quad w_t \sim N(0, \Sigma_w)$$

$$A_t \in \mathbb{R}^{n \times n} \quad B_t \in \mathbb{R}^{n \times d}$$

$$R^{(+)}(s_t, a_t) = -(S_t^T U_t S_t + a_t^T V_t a_t)$$

$$U_t \in \mathbb{R}^{n \times n} \quad V_t \in \mathbb{R}^{d \times d}$$

$$U_t \geq 0 \quad V_t \geq 0 \quad (\text{positive semidefinite, psd})$$

$$S_t^T U_t S_t \geq 0, \quad a_t^T V_t a_t \geq 0.$$

$$\Rightarrow R^{(+)}(s_t, a_t) \leq 0.$$

Assume the dynamic
is linear.

Helicopter: Want $S_t \approx 0$.

Choose $U_t = I$, $V_t = I$.

$$R(s_t, a_t) = -S_t^T S_t - a_t^T a_t$$
$$\approx \|S_t\|^2 - \|a_t\|^2$$

$$A = A_1 = A_2 = A_3 = \dots$$

pretend MDP is

$$B = B_1 = B_2 = B_3 = \dots$$

stationary for now.

$$S_{t+1} = A S_t + B a_t$$

$$s_0^{(i)} \xrightarrow{a_0^{(i)}} s_1^{(i)} \xrightarrow{a_1^{(i)}} s_2^{(i)} \xrightarrow{a_2^{(i)}} \dots \xrightarrow{a_T^{(i)}} s_T^{(i)} \quad (i=1, \dots, m)$$

$$\underset{A, B}{\operatorname{argmin}} \sum_{t=1}^T \|s_{t+1} - (As_t + Ba_t)\|^2.$$

Linearize non-linear model.

$$s_{t+1} = f(s_t, a_t)$$

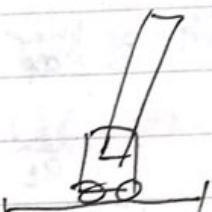
$$\begin{bmatrix} x_{t+1} \\ \dot{x}_{t+1} \\ \theta_{t+1} \\ \dot{\theta}_{t+1} \end{bmatrix} = f\left(\begin{bmatrix} x_t \\ \dot{x}_t \\ \theta_t \\ \dot{\theta}_t \end{bmatrix}, a\right)$$

$$s_{t+1} = f(s_t) \text{ pretend no action.}$$

$$s_{t+1} \approx f(\bar{s}_t) (\bar{s}_t - \bar{s}_t) + f(\bar{s}_t)$$

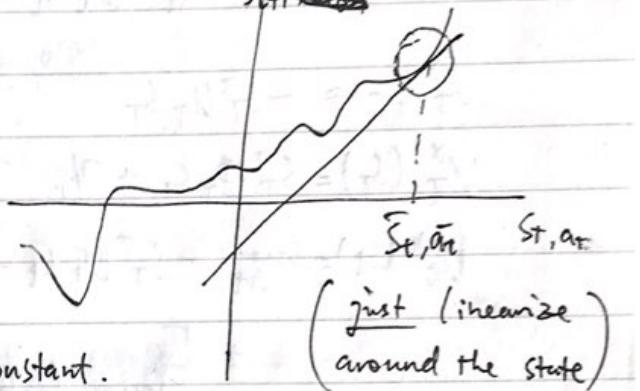
$$\begin{array}{l} \uparrow \\ s_{t+1} = As_t + Ba_t \\ \text{expected state near 0) } \\ \text{linearized near 0 state} \end{array}$$

$$\bar{s}_t = \text{constant.}$$



Download packets
to simulate IndP. on Internet

~~s_{t+1}~~



$$s_{t+1} = As_t + Ba_t + C \quad s = \begin{pmatrix} 1 \\ x \\ \dot{x} \\ \theta \\ \dot{\theta} \end{pmatrix}$$

$$\boxed{s_{t+1}} \approx f(\bar{s}_t, \bar{a}_t) + (\nabla_s f(\bar{s}_t, \bar{a}_t))^T ((\bar{s}_t - \bar{s}_t) + (\nabla_a f(\bar{s}_t, \bar{a}_t))^T (\bar{a}_t - \bar{a}_t))$$

\bar{s}_t, \bar{a}_t constant.

$$\text{Goal: } \max R^{(0)}(s_0, a_0) + R^{(1)}(s_1, a_1) + \dots + R^{(T)}(s_T, a_T)$$

$$\begin{aligned} V_T^*(s_T) &= \max_{a_T} R^{(T)}(s_T, a_T) \\ &= \max_{a_T} -s_T^T U_T s_T - a_T^T V_T a_T \end{aligned}$$

$$= \max_{a_T} -s_T^T U_T s_T \quad (\text{because } a_T^T V_T a_T \geq 0)$$

$$\tilde{U}_T^*(s_T) = \max_{a_T} R^{(T)}(s_T, a_T) = 0.$$

$a_T = 0 ?$

$$\underbrace{V_{t+1}^*}_{\text{DP}} \rightsquigarrow \underbrace{V_t^*}_{\text{DP}}$$

$$V_t^*(s_t) = \max_{a_t} R^{(t)}(s_t, a_t) + \sum_{s_{t+1}} P_{s_t a_t}(s_{t+1}) V_{t+1}^*(s_{t+1})$$

$$= \max_{a_t} R^{(t)}(s_t, a_t) + E_{s_{t+1} \sim P_{s_t a_t}} [V_{t+1}^*(s_{t+1})]$$

Suppose $\boxed{V_{t+1}^*(s_{t+1}) = S_{t+1}^\top \Phi_{t+1} s_{t+1} + \psi_{t+1}}$ | $\begin{array}{l} \Phi_{t+1} \in \mathbb{R}^{n \times n} \\ \psi_{t+1} \in \mathbb{R} \end{array}$

(can show) $V_t^*(s_t) = S_t^\top \Phi_t s_t + \psi_t$ for some Φ_t, ψ_t .

$$V_T^*(s_T) = -S_T^\top U_T S_T \quad \text{so } \Phi_T = -U_T \quad \psi_T = 0.$$

$$V_T^*(s_T) = S_T^\top \Phi_T s_T + \psi_T \quad (\text{Recursion...})$$

$$V_t^*(s_t) = \max_{a_t} -S_t^\top U_t s_t - A_t^\top V_t a_t.$$

$$+ E_{s_{t+1} \sim N(A_t s_t + B_t a_t, \Sigma_w)} [S_{t+1}^\top \underbrace{\Phi_{t+1} s_{t+1} + \psi_{t+1}}_{V_{t+1}^*(s_{t+1})}]$$

Simplifies to quadratic function of a_t
 $\Rightarrow a_t = \underbrace{(B_t^\top \Phi_{t+1} B_t - V_t)^{-1} B_t^\top \Phi_{t+1} A_t \cdot s_t}_{L_t}$

$$U_t^*(s_t) = \arg \max_{a_t} R^{(t)}(s_t, a_t) + E_{s_{t+1} \sim P_{s_t a_t}} [V_{t+1}^*(s_{t+1})]$$

$$= L_t \cdot s_t \quad (\text{NOT approx.})$$

→ plug back into D.P. recursion.

$$V_t^*(s_t) = S_t^\top \Phi_t s_t + \psi_t$$

$$\text{where } \Phi_t = A_t^\top (A_t - \Phi_{t+1} B_t (B_t^\top \Phi_{t+1} B_t - V_t)^{-1} B_t^\top \Phi_{t+1}) A_t - U_t$$

$$\psi_t = -\text{tr} \sum_w \Phi_{t+1} + \psi_{t+1}$$

Discrete time
Riccati equation