To summarize:

Initialize $\vec{\Phi}_T = -U_T$. $\gamma_T = 0$.

Recursive calculate $\vec{\Phi}_t$, $\gamma_t$ using $\vec{\Phi}_{t+1}$, $\gamma_{t+1}$ with D.T.R. $\vec{\Phi}_8^4$

$\qquad\qquad\qquad\qquad\qquad\qquad$ (for $t \geq T-1, T-2, \ldots, 0$)

Compute $L_t$ using $\vec{\Phi}_{t+1}$, $\gamma_{t+1}$.

$\qquad \pi_t^*(s_t) = L_t s_t$ $\qquad\qquad$ ( scale polynomial complexity

$\qquad V_t^*(s_t) = s_t^T \vec{\Phi}_t s_t + \gamma_t$. $\qquad$ $\quad$ to high dimension )

+ apply variation of this to helicopter

| Lecture 19 |

    - Debugging R.L. algorithm

    - LQR

      - Differential dynamics programming (DDP)

    - Kalman filter

    - Linear Quadratic Gaussian (LQG)

| "Advice for applying Machine Learning" [PDF] |

$$\max \; \mathbb{E}[R(s_0, a_0) + \cdots + R(s_T, a_T)].$$

DP: 
$$V_T^*(s) = \max_{a_T} R(s_T, a_T)$$

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

$$\pi_t^*(s) = \arg\max_a R(s,a) + \sum_{s'} P_{sa}^{(t)}(s') V_{t+1}^*(s')$$
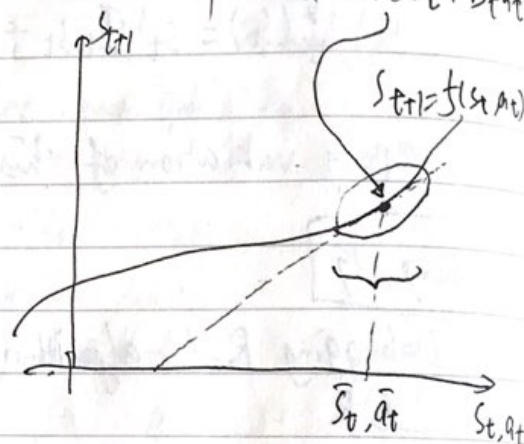
$$(V_T^*, V_{T-1}^*, V_{T-2}^*, \ldots, V_0^*)$$

LQR: $s_t \in \mathbb{R}^n$, $a_t \in \mathbb{R}^d$.

$$s_{t+1} = A_t s_t + B_t a_t + w_t \quad \leftarrow w_t \sim N(0, \Sigma_w)$$

$$\underbrace{\qquad}_{(\text{ignored})}$$

$$R(s_t, a_t) = -(s_t^T U_t s_t + a_t^T V_t a_t).$$

DP: 
$$V_t^*(s_t) = s_t^T \Phi_t s_t + \Psi_t \underbrace{\qquad}_{}$$

$$V_T^*: \quad \Phi_T = -Q_T \quad \Psi_T = 0$$

$$V_t^* \begin{cases} \Phi_t = -A_t\left(\Phi_{t+1} - \Phi_{t+1}B_t(B_t^T \Phi_{t+1}B_t - V_t)^{-1}B_t^T \Phi_{t+1}\right)A_t - Q_t \\ \Psi_t = -tr\,\Sigma_w \Phi_{t+1} + \Psi_{t+1} \quad (\text{ignored}) \end{cases}$$

→ Do depend on noise, larger noise cause worse value function.

$$\pi^*(s_t) = L_t s_t$$

$$L_t = (B_t^T \Phi_{t+1} B_t - V_t)^{-1} B_t^T \Phi_{t+1} A_t \; \leftarrow \text{Not depend on } \Psi$$

Special property of LQR systems
{
① You can forget $\Psi_t$'s !

② $\pi^*$'s don't depend on noise terms $w_t$'s.
}
↑
Kalman filter use this property.

(right column)

$$s_{t+1} = f(\bar{s}_t, \bar{a}_t)$$
$$+ (\nabla_s f(\bar{s}_t, \bar{a}_t))^T (s_t - \bar{s}_t)$$
$$+ (\nabla_a f(\bar{s}_t, \bar{a}_t))^T (a_t - \bar{a}_t)$$
↓
$$s_{t+1} = A_t s_t + B_t a_t$$

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



$\bar{s}_t, \bar{a}_t$    $s_t, a_t$

# Differential Dynamic Programming

$$S_{t+1} = f(S_t, a_t) \quad - \text{Simulator : nonlinear . deterministic}$$

## Have wanted trajectory

apply LQR on:
helicopter.
car.
chemical factory

(1) Come up with nominal trajectory
$$\bar{S}_0, \bar{a}_0, \bar{S}_1, \bar{a}_1, \ldots, \bar{S}_T, \bar{a}_T$$

(2) linearise $f$ around nominal Trajectory
i.e. $S_{t+1} \approx f(\bar{S}_t, \bar{a}_t) + (\nabla_s f(\bar{S}_t, \bar{a}_t))^T (S_t - \bar{S}_t)$
$$+ (\nabla_a f(\bar{S}_t, \bar{a}_t))^T (a_t - \bar{a}_t)$$

$$= A_t S_t + B_t a_t$$

$$\text{Expect}(S_t, a_t) \approx (\bar{S}_t, \bar{a}_t)$$

(3) Use LQR to get $\pi_t$

(4) Use simulator to get new nominal Trajectory
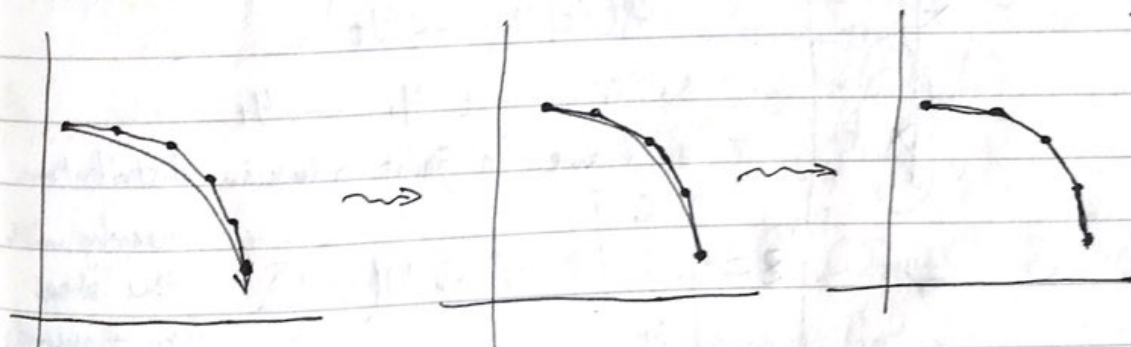i.e. $\bar{S}_0 = $ initial state
$$\bar{a}_t = \pi_t(\bar{S}_t)$$
$$\bar{S}_{t+1} = f(\bar{S}_t, \bar{a}_t)$$

$$\bar{S}_0, \bar{a}_0, \bar{S}_1, \bar{a}_1, \ldots, \bar{S}_T, \bar{a}_T$$

Linearize around new trajectory and repeat.



DDP:
local
optimal
search
algorithm.

This works well on the Ng's helicopter and works well on many problems.

# Kalman Filter & LQG.

Assume: <u>know</u> the state of system so far: $\bar{\pi}^k(s_t) = L_t s_t$

But when <u>Cannot</u> observe the state explicitly in some dynamic systems?

$$S_{t+1} = AS_t + W_t \qquad \text{(now forget control first)}$$

$$S_t = \begin{pmatrix} x_t \\ \dot{x}_t \\ \theta_t \\ \dot{\theta}_t \end{pmatrix} \qquad A = \begin{pmatrix} 1 & 1 & & \\ & 0.9 & & \\ & & 1 & 1 \\ & & & 0.9 \end{pmatrix} \begin{matrix} t \\ t \end{matrix}$$
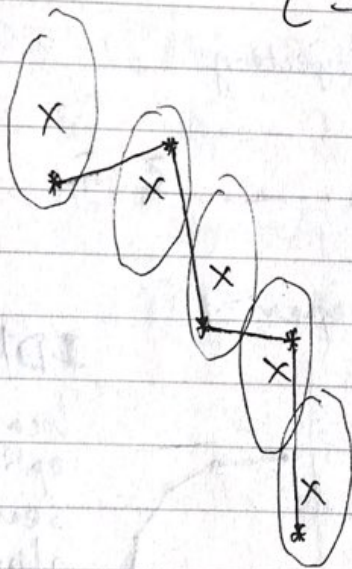
have a simulator and a radar to position the helicopter and estimate its states.

$$x_{t+1} = x_t + \dot{x}_t + \text{noise}$$
$$\dot{x}_{t+1} = 0.9 \dot{x}_t + \text{noise}$$

Observe

$$y_t = CS_t + V_t \qquad , \quad V_t \sim N(0, \Sigma_v)$$

$$C = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix}, \quad CS_t = \begin{bmatrix} x \\ y \end{bmatrix}.$$



* observation $y_t$,
X actual position.

estimate distribution on the state:

Want

$$P(s_t \mid y_1, \dots y_t)$$

$$s_0, s_1, \dots s_t, y_1, \dots, y_t$$

have a joint Gaussian distribution

$$Z = \begin{bmatrix} s_0 \\ s_1 \\ \vdots \\ s_t \\ y_1 \\ \vdots \\ y_t \end{bmatrix} \qquad Z \sim N(\mu, \Sigma)$$
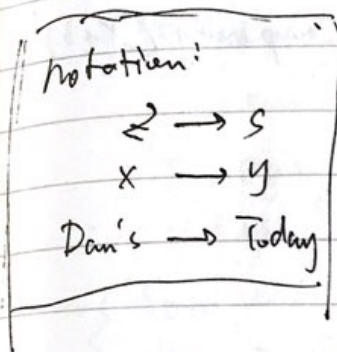
linearly with time steps (like: thousands)

Compute: $P(s_t \mid y_1, \dots y_t)$

Computationally inefficient.

radar

# Kalman Filter.

Dan's discussion on HMM's

Kalman Filter actually be a HMM.

Today: HMM with continous state rather than discrete states

$$P(s_t \mid y_1 \cdots y_t) \quad \Big\} \text{Predict}$$
$$P(s_{t+1} \mid y_1 \cdots y_t) \quad \Big\}$$
$$P(s_{t+1} \mid y_1 \cdots y_{t+1}) \quad \Big\} \text{Update}$$

## Predict step:

$$s_t \mid y_1 \cdots y_t \sim N(s_{t|t}, \Sigma_{t|t})$$

$s_t, y_t \longrightarrow$ true states/observations

Then
$$s_{t+1} \mid y_1 \cdots y_t \sim N(s_{t+1|t}, \Sigma_{t+1|t})$$

Where
$$s_{t+1|t} = A\, s_{t|t}$$
$$\Sigma_{t+1|t} = A\, \Sigma_{t|t} A^T + \Sigma_v$$

$s_{t|t} \quad s_{t+1|t}$

$\Sigma_{t|t} \quad \Sigma_{t+1|t}$

$\longrightarrow$ Computations

## Update step:

Get $s_{t+1} \mid y_1 \cdots y_{t+1} \sim N(s_{t+1|t+1}, \Sigma_{t+1|t+1})$

where $s_{t+1|t+1} = s_{t+1|t} + K_{t+1} \cdot (y_{t+1} - C\, s_{t+1|t})$

$$K_{t+1} = \Sigma_{t+1|t}\, C^T \left( C\, \Sigma_{t+1|t}\, C^T + \Sigma_v \right)^{-1}$$

$$\Sigma_{t+1|t+1} = \Sigma_{t+1|t} - \Sigma_{t+1|t}\, C^T \left( C\, \Sigma_{t+1|t}\, C^T + \Sigma_v \right)^{-1} \cdot C \cdot \Sigma_{t+1|t}$$

$s_{t+1|t+1}$ is our "best" estimate for $s_{t+1}$.

$$Z = \begin{bmatrix} s_0 \\ \vdots \\ s_t \\ y_1 \\ \vdots \\ y_t \end{bmatrix} \qquad Z \sim N(\mu, \Sigma)$$

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

If Compute marginal dist. : Time complexity : $O(t^3)$

KF : $O(1)$ for every step.



$$P(s_1|y_1) \rightsquigarrow P(s_2|y_1, y_2) \rightsquigarrow \boxed{P(s_3|y_1, y_2, y_3)}$$

<u>Putting these together ( KF + LQR = LQG ).</u>

$$s_{t+1} = A s_t + B a_t + w_t \qquad w_t \sim N(0, \Sigma_w)$$
$$y_t = C s_t + v_t \qquad v_t \sim N(0, \Sigma_v)$$

Use KF to estimate state

$$s_{0|0} = s_0 \qquad \Sigma_{0|0} = 0 \qquad \text{for} \quad s_0 \sim N(s_{0|0}, \Sigma_{0|0}).$$

predict $\begin{cases} s_{t+1|t} = A s_{t|t} + \boxed{B a_t} \\ \Sigma_{t+1|t} = A \Sigma_{t|t} A^T + \Sigma_v \end{cases}$

Compute $L_t$'s using LQR   (Assuming observed states)

$$a_t = L_t s_t \longrightarrow a_t = L_t \cdot s_{t|t}$$

" $s_t = s_{t|t} + noise$ "   This is optimal. Due to <u>seperate</u> principle.

Only hold true for spec. case like LQG.

(estimate state and directly plug-in estimation into controller)

For many other systems (nonlinear or other changes in LQG) this is NOT hold true. i.e. not optimal if you do this.

# Lecture 20 (The last lecture of CS229)

- POMDPs (Partially observed MDPs)
- Policy search
- Reinforce algorithm
- Pegasus algorithm
- Conclusion

$$\begin{cases} S_{t+1} = A S_t + B a_t + w_t \\ Y_t = C S_t + v_t \qquad \leftarrow \text{observation} \end{cases}$$

Actions $\quad a_t = L_t S_t$

Compute $\quad S_{t|t}$ (estimate for $s_t$)

Kalman filter: $\quad S_t | y_1 \cdots y_t \sim N(S_{t|t}, \Sigma_{t|t})$.

Actions: $\quad a_t = L_t S_{t|t}$

---

\* Find the optimal policy for POMDP is NP-hard.

<u>POMDP</u>: $(S, A, Y, \{P_{sa}\}, \{O_s\}, T, R)$

$\quad Y$ - set of possible observations

$\quad O_s$ - observation distributions.

At each step, observe $Y_t \sim O_{S_t}$ (if in state $S_t$)

---

<u>Policy Search</u> (Direct policy search)

One of most effective algo's for Full-Observed MDPs and Partial-Observed MDPs (POMDPs).

Define a set $\overline{\Pi}$ of policies, Search for a goal $\pi \in \overline{\Pi}$.

(c.f. Define a set $\mathcal{H}$ of hypothesis, search for a good $h \in \mathcal{H}$)

New definition: A stochastic policy is a function.
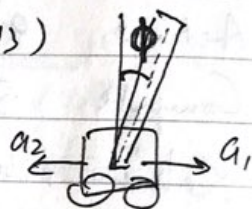
$\pi : S \times A \mapsto \mathbb{R}$ when $\pi(s,a)$ is probability of taking action "a" in state s.

$\left( \sum_a \pi(s,a) = 1 . \quad \pi(s,a) \geq 0 \right)$

---

execute $\pi$:     In state s. take action $a_1 : a_2 : a_3$

$= \pi(s,a_1) : \pi(s,a_2) : \pi(s,a_3)$
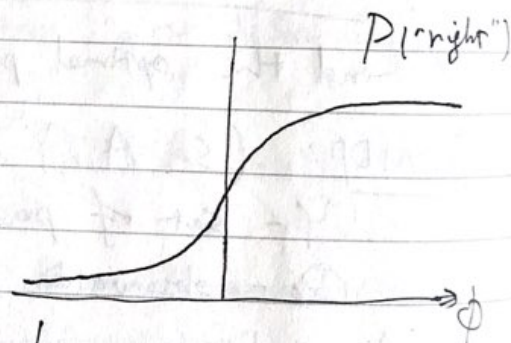
$\pi_\theta(s,a_1) = \dfrac{1}{1+e^{-\theta^T s}}$

$\pi_\theta(s,a_2) = 1 - \dfrac{1}{1+e^{-\theta^T s}}$

E.g. $S = \begin{bmatrix} 1 \\ x \\ \dot{x} \\ \phi \\ \dot{\phi} \end{bmatrix}$    $\theta = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$

$P("a = right") = \dfrac{1}{1+e^{-\theta^T s}} = \dfrac{1}{1+e^{-\phi}}$

P("right")

Goal: $\max_\theta E\left[ R(s_0, a_0) + \cdots + R(s_T, a_T) \mid \pi_\theta . s_0 \right]$

$\theta_1 \cdots \theta_d : \quad \pi_\theta(s, a_i) = \dfrac{e^{\theta_i^T s}}{\sum_{j=1}^{d} e^{\theta_j^T s}}$    (softmax)

# Reinforce Algorithm (isn't exactly the reinforce algorithm, but as originally presented by Ron Williams, but it captures its essence).

Assume $s_0$ is some fixed initial state,

$$\max E[R(s_0, a_0) + \cdots + R(s_T, a_T)]$$

$$= \sum_{\substack{s_0 a_0 \cdots \\ s_T a_T}} P(s_0, a_0, s_1, a_1 \cdots s_T a_T) \left[ R(s_0, a_0) + \cdots + R(s_T, a_T) \right].$$

$$= \sum_{\substack{s_0 a_0 \cdots \\ s_T a_T}} P(s_0) \cdot \pi_\theta(s_0, a_0) P_{s_0 a_0}(s_1) \pi_\theta(s_1, a_1) \cdots \cdots \pi_\theta(s_T, a_T)$$

$$* \underbrace{\left[ R(s_0, a_0) + \cdots + R(s_T, a_T) \right]}_{\text{Payoff}}.$$

Loop: {

Sample $s_0, a_0, s_1, a_1, \cdots, s_T, a_T$

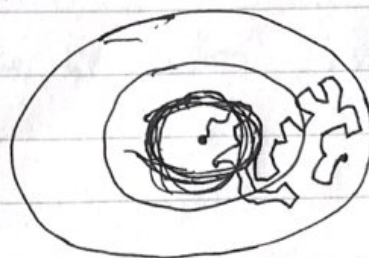Compute payoff $= R(s_0, a_0) + \cdots + R(s_T, a_T)$

Update:

$$\theta := \theta + \alpha \left[ \frac{\nabla_\theta \pi_\theta(s_0, a_0)}{\pi_\theta(s_0, a_0)} + \cdots + \frac{\nabla_\theta \pi_\theta(s_T, a_T)}{\pi_\theta(s_T, a_T)} \right] \cdot \text{payoff}$$

}

$$E\left[ \left( \frac{\nabla_\theta \pi_\theta(s_0, a_0)}{\pi_\theta(s_0, a_0)} + \cdots + \frac{\nabla_\theta \pi_\theta(s_T, a_T)}{\pi_\theta(s_T, a_T)} \right) \cdot \text{payoff} \right].$$



converge slowly ...  → Stochastic gradient ascent algorithm

$$\nabla_\theta E[\text{payoff}]$$

$$= \sum_{\substack{s_0 a_0 \cdots \\ s_T a_T}} \left[ P(s_0) \big( \nabla_\theta \pi_\theta(s_0, a_0) \big) P_{s_0 a_0}(s_1) \pi_\theta(s_1 a_1) \cdots \pi_\theta(s_T a_T) \right.$$

$$+ P(s_0) \pi_\theta(s_0, a_0) P_{s_0 a_0}(s_1) \big( \nabla_\theta \pi_\theta(s_1 a_1) \big) \cdots \pi_\theta(s_T a_T)$$

$$+ \cdots$$

$$\left. + P(s_0) \pi_\theta(s_0, a_0) P_{s_0 a_0}(s_1) \pi_\theta(s_1, a_1) \cdots \big( \nabla_\theta \pi_\theta(s_T, a_T) \big) \right]$$

$$* \text{payoff}$$

$$\frac{d}{d\theta} f(\theta) g(\theta) h(\theta) = f'(\theta) g(\theta) h(\theta) + f(\theta) g'(\theta) h(\theta) + f(\theta) g(\theta) h'(\theta)$$

4 steps derivation

$$= \sum_{\substack{s_0 a_0 \cdots \\ s_T a_T}} P(s_0) \pi_\theta(s_0, a_0) P_{s_0 a_0}(s_1) \pi_\theta(s_1, a_1) \cdots \cdots \pi_\theta(s_T a_T)$$

$$\cdot \left[ \frac{\nabla_\theta \pi_\theta(s_0, a_0)}{\pi_\theta(s_0, a_0)} + \frac{\nabla_\theta \pi_\theta(s_1, a_1)}{\pi_\theta(s_1, a_1)} + \cdots + \frac{\nabla_\theta \pi_\theta(s_T, a_T)}{\pi_\theta(s_T, a_T)} \right] \cdot \text{payoff}$$

$$= \sum_{\substack{s_0 a_0 \cdots \\ s_T a_T}} P(s_0, a_0, s_1, a_1 \cdots s_T a_T) \left[ \frac{\nabla_\theta \pi_\theta(s_0, a_0)}{\pi_\theta(s_0, a_0)} + \cdots + \frac{\nabla_\theta \pi_\theta(s_T, a_T)}{\pi_\theta(s_T, a_T)} \right] \cdot \text{payoff}$$

$$= E\left[ \left( \frac{\nabla_\theta \pi_\theta(s_0, a_0)}{\pi_\theta(s_0, a_0)} + \cdots + \frac{\nabla_\theta \pi_\theta(s_T, a_T)}{\pi_\theta(a_T, a_T)} \right) \cdot \text{payoff} \right]$$

① policy search algo.'s are esp. effective when you can chose a simple policy $\pi$: <u>For the proble if exist a simple function ($LP/Logit R$) that maps features of states to the action)</u>.

    ↳ a proper policy class $\underline{\pi}$.     eg: Inv. P. ~~Flying a helicopter~~

                       eg: <u>Low Level control/reflexes</u>:

                               Flying a helicopter. driving a car.

② If the problem requires <u>long multi-step reasoning</u> (eg: game of chess), <u>high level decision (less instinctual) making</u>. Use value function approximation approaches instead

Use approximation $\hat{s}$ of s. (Could be $\hat{s} = S_{t/t}$ from KF).

$$\pi_\theta(s,a) = \frac{1}{1 + e^{-\theta^T \hat{s}}}$$

                       (can use policy search algo's on POMDP)

                       (often reasonably effective to POMDPs).

* Reinforce algo. often works well, but is often ~~extremely~~ <u>slow</u>.

                               (because of <u>noise</u>)

         such as <u>1-100 million iterations</u>.

* So : Sample $s_0, a_0, s_1, a_1, \ldots, s_T, a_T$.

         10 million times: always on a simulator, not a physical device like robot.

# Pegasus policy search

[ Policy search : Pegasus     [PPT] ]

Actually use on Ng's autonomous helicopter flight for many years.

<u>Pegasus</u>: Policy Evaluation of Gradient And Search Using Scenarios.

<u>Scenarios</u>: Fixed random numbers (random seq's).
generated from random number generator.

* Main idea: Evaluate many times/scenarios on policy,
(for opt. a deterministic function), then average.

.* Scale well even to fairly large problems. (high-dim. statespace)

* Key in RL: sequential decision making
(consider long-term consequences).

Other e.g:

medical decision making :   seq. of treatments
[queues] { bank         :   multiple queues (wait time loss)
          { assembly line :   objects in queues
financial decision making:  sell off stocks.
[OR problems] factory automation:  opt. throughput/cost.

<u>R.L. Applications</u>

Little Dog robot : by Ziko Coulter[TA] & Peter Abiel[PhD].
(run similar to: approx. value func.)
Legged wheeled robot: wheeled robots: very fuel-efficient.
(cars. trucks)
by Lockheed Martin Cooperation.
Helicopter: by Peter Abiel & Adam Coates.

Machine Learning widely used in:

Industry management.
Optimize computer architecture.
Network security.
Robotics
Computer vision
Computational biology.
Aerospace
Natural Language Understanding (NLP)..:

## Choose AI classes

Stanford has one of the best & broadest sets of AI classes:
Learn more about AI, other fields which often apply learning
algorithms to problems.

- CS221. Overview of AI (by Ng)
✗ CS228. Probabilistic models in AI (by Daphne Koller).
        closest in spirit to 229.
                HIGHLY RECOMMENDED (so as to Ng's PhDs).

- EE366. Convex Optimization. (by Stephen Boyd).
- CS294. Project Course. (by Ng)

[END]