CS-GY 6923: Lecture 14 Finish Semantic Embeddings, Modern Image Generation, Reinforcement Learning

NYU Tandon School of Engineering, Prof. Christopher Musco

SEMANTIC EMBEDDING

Goal: Learn mapping from inputs to numerical vectors such that similar inputs map to similar vectors (e.g., with high inner product).



Goal: Learn mapping from inputs to numerical vectors such that similar inputs map to similar vectors (e.g., with high inner product).



For words, the mapping is typically just a lookup table.

HOW TO GET EMBEDDINGS?



For documents or words, earliest approaches were based on latent semantic analysis (PCA on term document matrix).

More modern word embedding recipe:

- Choose similarity metric k(word_i, word_j) which can be computed for any pair of words.
- 2. Construct similarity matrix $\mathbf{M} \in \mathbb{R}^{n \times n}$ with $\mathbf{M}_{i,j} = k(word_i, word_j)$.
- 3. Find low rank approximation $\mathbf{M} \approx \mathbf{Y}^T \mathbf{Y}$ where $\mathbf{Y} \in \mathbb{R}^{k \times n}$.
- 4. Columns of Y are word embedding vectors.

We expect that $\langle \mathbf{y}_i, \mathbf{y}_j \rangle$ will be larger for more similar words.

Common choice for similarity metric is to use co-occurence frequency in windows.

The girl walks to her dog to the park. It can take a long time to park your car in NYC. The dog park is always crowded on Saturdays.

The girl walks to her dog to the park. It can take a long time to park your car in NYC. The dog<mark>park is always crowded</mark>on Saturdays.

The girl walks to her dog to the park. It can take a long time to park your car in NYC. The dog park is always crowded on Saturdays.

	dog	park	crowded	the
gop	0	2	0	3
park	2	0	1	2
crowded	0	1	0	0
the	3	2	0	0

Usually followed by some tranformation or normalization. E.g., $k(word_i, word_j) = \frac{p(i,j)}{p(i)p(j)}$.

Current state of the art models: GloVE, word2vec.

- **word2vec** was originally presented as a shallow neural network model, but it is equivalent to matrix factorization method (Levy, Goldberg 2014).
- For word2vec, similarity metric is the "point-wise mutual information": $\log \frac{p(i,j)}{p(i)p(j)}$.

Common to use <u>pre-trained</u> word vectors:

 Compilation of many sources: https://github.com/3Top/word2vec-api

CAVEAT ABOUT FACTORIZATION



SVD will not return a symmetric factorization in general. In fact, if **M** is not positive semidefinite¹ then the optimal low-rank approximation does not have this form.

¹I.e., *k*(*word*_{*i*}) is not a positive semidefinite kernel.

CAVEAT ABOUT FACTORIZATION



- For each word *i* we get a left and right embedding vector
 w_i and y_i. It's reasonable to just use one or the other.
- If (y_i, y_j) is large and positive, we expect that y_i and y_j have similar similarity scores with other words, so they typically are still related words.
- Another option is to use as your features for a word the concatenation [w_i, y_i]

The same approach used for word embeddings can be used to obtain meaningful numerical features for any other data where there is a natural notion of similarity.



For example, the items could be nodes in a social network graph. Maybe be want to predict an individuals age, level of interest in a particular topic, political leaning, etc.

NODE EMBEDDINGS



Generate random walks (e.g. "sentences" of nodes) and measure similarity by node co-occurence frequency.



Again typically normalized and apply a non-linearity (e.g. log) as in word embeddings.



Popular implementations: **DeepWalk**, **Node2Vec**. Again initially derived as simple neural network models, but are equivalent to matrix-factorization (Qiu et al. 2018).

We can also create embeddings that represent different types of data. OpenAI's clip architecture:



Goal: Train embedding architectures so that $\langle T_i, I_j \rangle$ are similar if image and sentence are similar.

CLIP TRAINING

What do we use as ground truth similarities during training? Sample a batch of sentence/image pairs and just use identity matrix.



My new puppy!	1	0	0
Best dim sum ever.	0	1	0
NYC in the rain.	0	0	1

This is called <u>contrastive learning</u>. Train unmatched text/image pairs to have nearly orthogonal embedding vectors.

CLIP FOR ZERO-SHOT LEARNING

Learning Transferable Visual Models From Natural Language Supervision

Alec Radford ^{*1} Jong Wook Kim ^{*1} Chris Hallacy ¹ Aditya Ramesh ¹ Gabriel Goh ¹ Sandhini Agarwal ¹ Girish Sastry ¹ Amanda Askell ¹ Pamela Mishkin ¹ Jack Clark ¹ Gretchen Krueger ¹ Ilya Sutskever ¹



2021 result: 76% accuracy on ImageNet image classification challenge with <u>no labeled training data.</u>

IMAGE SYNTHESIS (TEASER)

AUTOENCODERS LEARN COMPRESSED REPRESENTATIONS



 $f(\mathbf{x}) = d(e(\mathbf{x}))$ projects an image **x** closer to the space of natural images.

Suppose we want to generate a random natural image. How might we do that?

• **Option 1**: Draw each pixel value in **x** uniformly at random. Draws a random image from *A*.



• **Option 2**: Draw **x** randomly from *S*, the space of images representable by the autoencoder.



How do we randomly select an image from \mathcal{S} ?

Autoencoder approach to generative ML: Feed random inputs into decode to produce random realistic outputs.



Main issue: most random inputs words will "miss" and produce garbage results.

AUTOENCODERS FOR DATA GENERATION



Variational auto-encoders attempt to resolve this issue.

Developed from a different perspective than regular autoencoders. Make the data generation goal more explicit.

- Train a neural network G_{θ} that takes in a length k code word, **z**, and outputs an image.
- + Assume $\boldsymbol{z} \sim \mathcal{N}(\boldsymbol{0},\boldsymbol{I}).$ I.e., a random Gaussian vector.
- Goal is to maximize probability of producing a "natural image".

First attempt: Given training data $\mathbf{x}_1, \ldots, \mathbf{x}_n$,

$$\max_{\theta} \int \mathbb{1}[G_{\theta}(\mathbf{z}) = \mathbf{x}_{i} \text{ for some } i] \cdot p(\mathbf{z})d\mathbf{z}$$
$$= \max_{\theta} \mathbb{E}_{\mathbf{z}} \mathbb{1}[G_{\theta}(\mathbf{z}) = \mathbf{x}_{i} \text{ for some } i]$$

Issues: Super brittle, impossible to train.

Bayesian approach: assume each \mathbf{x}_i is of the form $G_{\theta}(\mathbf{z}) + \sigma \mathcal{N}(0, \mathbf{I})$ for randomly chosen \mathbf{z} . Choose parameters, θ , to maximize the likelihood of the data:

$$\begin{aligned} \max_{\theta} \prod_{i=1}^{n} p(\mathbf{x}_{i}) &= \max_{\theta} \prod_{i=1}^{n} \int p(\mathbf{x}_{i} \mid \mathbf{z}) \cdot p(\mathbf{z}) d\mathbf{z} \\ &= \max_{\theta} \sum_{i=1}^{n} \log \int p(\mathbf{x}_{i} \mid \mathbf{z}) \cdot p(\mathbf{z}) d\mathbf{z} \\ &= \min_{\theta} \sum_{i=1}^{n} -\log \int e^{-\|\mathbf{x}_{i} - G_{\theta}(\mathbf{z})\|_{2}^{2}/2\sigma^{2}} \cdot p(\mathbf{z}) d\mathbf{z} \end{aligned}$$

VARIATIONAL AUTOENCODERS

$$\max_{\theta} \prod_{i=1}^{n} p(\mathbf{x}_{i}) = \max_{\theta} \prod_{i=1}^{n} \int p(\mathbf{x}_{i} \mid \mathbf{z}) \cdot p(\mathbf{z}) d\mathbf{z}$$
$$= \max_{\theta} \sum_{i=1}^{n} \log \int p(\mathbf{x}_{i} \mid \mathbf{z}) \cdot p(\mathbf{z}) d\mathbf{z}$$

How to deal with the integral? Very common approach in generative modeling (beyond VAEs): <u>Monte Carlo</u> <u>approximation</u>. Draw samples z_1, \ldots, z_m and observe that:

$$pprox \max_{\theta} \sum_{i=1}^{n} \log \left(\frac{1}{m} \sum_{i=1}^{m} p(\mathbf{x}_i \mid \mathbf{z}_i) \right).$$

This approach does not work out of the box. The issue is that the integral will be very poorly approximated by sampling:

$$\int p(\mathbf{x}_i \mid \mathbf{z}) \cdot p(\mathbf{z}) d\mathbf{z} \quad \not\approx \quad \sum_{i=1}^m p(\mathbf{x}_i \mid \mathbf{z}_i).$$

Second key idea: Importance sampling. For any distribution q(z), $p(x_i) = \int p(x_i \mid z) \cdot p(z) dz = \int q(z) \frac{p(x_i \mid z)}{q(z)} \cdot p(z) dz$

Draw $\mathbf{z}_1, \ldots, \mathbf{z}_m$ from $q(\mathbf{z})$ and estimate:

$$p(\mathbf{x}_i) \approx \frac{1}{m} \sum_{i=1}^m \frac{p(\mathbf{x}_i \mid \mathbf{z})}{q(\mathbf{z})} \cdot p(\mathbf{z}).$$

We can choose a different distribution for each \mathbf{x}_i . I.e., choose q_1, \ldots, q_n . Ideally, want q_i to be higher for \mathbf{z} that are more likely to generate \mathbf{x}_i . Ideal choice is $q_i(\mathbf{z}) = p(\mathbf{z} \mid \mathbf{x}_i)$.

$$\frac{1}{m}\sum_{i=1}^{m}\frac{p(\mathbf{x}_i \mid \mathbf{z})}{q_i(\mathbf{z})} \cdot p(\mathbf{z}) = \frac{1}{m}\sum_{i=1}^{m}\frac{p(\mathbf{x}_i \mid \mathbf{z}) \cdot p(\mathbf{z})}{p(\mathbf{z} \mid \mathbf{x}_i)}$$

Typical VAE approach: Assume q_i is parameterized as a multivariate Gaussian distribution with mean $\mu_i \in \mathbb{R}^k$ and variances $\boldsymbol{\Sigma} = [\sigma_1^2, \dots, \sigma_k^2]$. Train a model (e.g., neural network) that maps \mathbf{x}_i to $\mu_i, \boldsymbol{\Sigma}_i$.

Simulateously minimize distance between q_i and $p(\mathbf{z} | \mathbf{x}_i)$ (typically using KL divergence) and maximize $\sum_{i=1}^{n} p(\mathbf{x}_i)$, where $p(\mathbf{x}_i)$ is approximated via importance sampling.

Lots of details here! Link to some good notes by Brian Kang.

VAEs are not really autoencoders. Not designed to map an input **x** to an approximation **x**̃. But, their final architecture ends up resembling that of an autoencoder:



VAE's give very good results, but tend to produce images with immediately recognizable flaws (e.g. soft edges, high-frequency artifacts).





Lots of efforts to hand-design regularizers that penalize images that don't look realisitic to the human eye.

Main idea behind GANs: Use machine learning to automatically encourage realistic looking images.

 $\min_{\theta} L(\theta) + P(\theta)$

GENERATIVE ADVERSARIAL NETWORKS (GANS)



Let $\mathbf{x}_1, \ldots, \mathbf{x}_n$ be real images and let $\mathbf{z}_1, \ldots, \mathbf{z}_m$ be random code vectors. The goal of the discriminator is to output a number between [0, 1] which is close to 0 if the image is fake, close to 1 if it's real.

Train weights of discriminator D_{θ} to minimize:

$$\min_{\boldsymbol{\theta}} \sum_{i=1}^{n} -\log \left(D_{\boldsymbol{\theta}}(\mathbf{x}_{i}) \right) + \sum_{i=1}^{m} -\log \left(1 - D_{\boldsymbol{\theta}}(G_{\boldsymbol{\theta}'}(\mathbf{z}_{i})) \right)$$
²⁹

GENERATIVE ADVERSARIAL NETWORKS (GANS)



Goal of the generator $G_{\theta'}$ is the opposite. We want to maximize:

$$\max_{\boldsymbol{\theta}'} \sum_{i=1}^{m} -\log\left(1 - D_{\boldsymbol{\theta}}(G_{\boldsymbol{\theta}'}(\mathbf{z}_i))\right)$$

This is called an "adversarial loss function". *D* is playing the role of the adversary.

$$\boldsymbol{\theta}^*, \boldsymbol{\theta}'^* \text{ solve } \min_{\boldsymbol{\theta}} \max_{\boldsymbol{\theta}'} \sum_{i=1}^n -\log \left(D_{\boldsymbol{\theta}}(\mathbf{x}_i) \right) + \sum_{i=1}^m -\log \left(1 - D_{\boldsymbol{\theta}}(G_{\boldsymbol{\theta}'}(\mathbf{z}_i)) \right)$$

This is called a minimax optimization problem. <u>Really tricky to</u> solve in practice.

- **Repeatedly play:** Fix one of *θ* or *θ'*, train the other to convergence, repeat.
- Simultaneous gradient descent: Run a single gradient descent step for each of θ, θ' and update D and G accordingly. Difficult to balance learning rates.

State of the art until a few years ago.



DIFFUSION

Auto-encoder/GAN approach: Input noise, map directly to image.

Diffusion: Slowly move from noise to image.



We will post a demo for generating MNIST digits via diffusion.





Tons of other work going on in image generation. One key topic is "class conditioned" generation:

2222222222 3333 222222222 3333333333333 *333333*3 2 2 3 3 7 7 7 7 7 7 3 3 3 3 3
Can also condition on another image...



SEMANTIC EMBEDDINGS + GENERATIVE MODELS

Or a sentence...



"A chair that looks like an avocado"

SEMANTIC EMBEDDINGS + GENERATIVE MODELS

Or a sentence...



"A diagram that explains variational autoencoders"

REINFORCEMENT LEARNING (TEASER)

Rest of lecture: Give flavor of the area and insight into <u>one</u> algorithm (Q-learning) which has been successful in recent years.

Basic setup:

- Agent interacts with environment over time 1, ..., t.
- Takes repeated sequence of **actions**, a_1, \ldots, a_t which effect the environment.
- **State** of the environment over time denoted s_1, \ldots, s_t .
- Earn **rewards** r_1, \ldots, r_t depending on actions taken and states reached.
- Goal is to maximize reward over time.

REINFORCEMENT LEARNING EXAMPLES

Classic inverted pendulum problem:



• Agent: Cart/software controlling cart.

• **State:** Position of the car, pendulum head, etc.

- Actions: Move cart left or move right.
- **Reward:** 1 for every time step that $|\theta| < 90^{\circ}$ (pendulum is upright). 0 when $|\theta| = 90^{\circ}$

This problem has a long history in **Control Theory.** Other applications of classical control:

- · Semi-autonomous vehicles (airplanes, helicopters, drones, etc.)
- Industrial processes (e.g. controlling large chemical reactions)
- Robotics



control theory : reinforcement learning :: stats : machine learning

REINFORCEMENT LEARNING EXAMPLES

Strategy games, like Go:



- **State:** Position of all pieces on board.
- Actions: Place new piece.

• **Reward:** 1 if in winning position at time *t*. 0 otherwise.

This is a <u>sparse reward problem</u>. Payoff only comes after many times steps, which makes the problem very challenging.

REINFORCEMENT LEARNING EXAMPLES

Video games, like classic Atari games:



- State: Raw pixels on the screen (sometimes there is also hidden state which can't be observed by the player).
- Actions: Actuate controller (up,down,left,right,click).
- **Reward:** 1 if point scored at time *t*.

Model problem as a Markov Decision Process (MDP):

- S : Set of all possible states. # of states is |S|.
- $\cdot \ \mathcal{A}$: Set of all possible actions. # of actions is $|\mathcal{A}|.$
- + $\mathcal R$: Set of possible rewards. Could have $\mathcal R=\mathbb R.$
- Reward function

 $R(s, a) : S \times A \rightarrow \text{ probability distribution over } \mathcal{R}. r_t \sim R(s_t, a_t).$

 \cdot State transition function

 $P(s, a) : S \times A \rightarrow \text{ probability distribution over } S. s_{t+1} \sim P(s_t, a_t).$

Why is this called a <u>Markov</u> decision process? What does the term Markov refer to?

MATHEMATICAL FRAMEWORK FOR RL

Goal: Find a **policy** $\Pi : S \to A$ from states to actions which maximize expected cumulative reward.

- Start is state s₀.
- For $t = 0 \dots, T$
 - $r_t \sim R(s_t, \Pi(s_t))$.
 - $s_{t+1} \sim P(s_t, \Pi(s_t))$.

The **time horizon** *T* could be short (game with fixed number of steps), very long (stock investing), or infinite. Goal is to maximize:

$$reward(\Pi) = \mathbb{E}\sum_{t=0}^{T} r_t$$

 $[s_0, a_0, r_0], [s_1, a_1, r_1], \dots, [s_t, a_t, r_t]$ is called a **trajectory** of the MDP under policy Π .²

 2 It turns out that it is always optimal to use a fixed policy. There is no benefit to changing Π over time. We will discuss this shortly.

FLEXIBILITY OF MDPS

- Can be used to model time-varying environments. Just add time *t* to the state vector.
- Can be used to model games where actions have different effect if play in sequence (e.g. combo in a video game).
 Just add list of previous few actions to state.
- Can be used to model two-player games. Model adversary as part of the transition function.



SIMPLE EXAMPLE: GRIDWORLD



- $r_t = -.01$ if not at an end position. ± 1 if at end position.
- $P(s_t, a)$: 70% of the time move in the direction indicated by *a*. 30% of the time move in a random direction.

What is the optimal policy Π ?

SIMPLE EXAMPLE: GRIDWORLD



- $r_t = -.5$ if not at an end position. ± 1 if at end position.
- $P(s_t, a)$: 70% of the time move in the direction indicated by *a*. 30% of the time move in a random direction.

What is the optimal policy Π ?

For infinite or very long times horizon games (large T), we often introduce a **discount factor** γ and seek instead to take actions which minimize:



where $r_t \sim R(s_t, \Pi(s_t))$ and $s_{t+1} \sim P(s_t, \Pi(s_t))$ as before.

 $\gamma \rightarrow$ 1: No discount. Standard MDP expected reward.

 $\gamma \rightarrow$ 0: Care about short term reward more.

From now on assume $T = \infty$. We can do this without loss of generality by adding a time parameter to state and moving into an "end state" with no additional rewards once the time hits *T*.

Value function: Measures the expected return if we start in state s and follow policy Π .

$$V^{\Pi}(s) = \mathbb{E}_{\Pi, s_0 = s} \sum_{t \ge 0} \gamma^t r_t$$

Let $\Pi_s^* = \arg \max V^{\Pi}(s)$. If we are in state *s*, <u>at any point</u>, we should always take action $\Pi_s^*(s)$.

Value function:

$$V^{\Pi}(s) = \mathbb{E}_{\Pi, s_0 = s} \sum_{t \ge 0} \gamma^t r_t$$

Claim: Let $\Pi_s^* = \arg \max V^{\Pi}(s)$. If we are in state *s*, <u>at any point</u>, we should always take action $\Pi_s^*(s)$.

Proof: Suppose we has already taken j - 1 steps and seen trajectory $[s_0, a_0, r_0], \ldots, [s_j, a_j, r_j]$. Then our expected reward is:

$$r_{0} + \gamma r_{1} + \ldots + \gamma^{j-1} r_{j-1} + \mathbb{E}_{\Pi} \sum_{t \ge j} \gamma^{t} r_{j}$$

= $r_{0} + \gamma r_{1} + \ldots + \gamma^{j-1} r_{j-1} + \gamma^{j} \cdot \mathbb{E}_{\Pi} \sum_{t \ge 0} \gamma^{t} r_{t+j}$
= $r_{0} + \gamma r_{1} + \ldots + \gamma^{j} r_{j} + \gamma^{j} \cdot V^{\Pi}(s_{j})$

Value function:

$$V^{\Pi}(s) = \mathbb{E}_{\Pi, s_0 = s} \sum_{t \ge 0} \gamma^t r_t$$

Claim: Let $\Pi_s^* = \arg \max V^{\Pi}(s)$. If we are in state s, <u>at any point</u>, we should always take action $\Pi_s^*(s)$.

Consequence: there is a <u>single</u> optimal policy Π^* which simultaneously maximizes $V^{\Pi}(s)$ for all s. I.e. $\Pi_1^* = \Pi_2^* = \ldots = \Pi_{|S|}^* = \Pi^*$. We do not need to change the policy over time to maximize expected reward.

Goal in RL is to find this optimal policy $\Pi^*.$

Full information: We know *S*, *A*, the transition function *P* and reward function *R*. Sometimes called the "planning" problem.

Reinforcement Learning setting: We do not know *P* or *R*, but we can repeatedly play the MDP, running whatever policy we like.

Let $V^*(s) = V^{\Pi^*}(s)$. This function is equal to the expected future reward if we play <u>optimally</u> starting in state *s*.



VALUE ITERATION

In the <u>full information</u> setting, if we knew V^* we can easily find the optimal policy Π :



$$\Pi^*(s) = \arg\max_{a} \sum_{s',r} \cdot \Pr(s',r \mid s,a)[r + \gamma V^*(s')]$$

*V**(*s*) satisfies what is called a <u>Bellman equation</u>:

$$V^*(s) = \max_{a} \sum_{s',r} \cdot \Pr(s',r \mid s,a)[r + \gamma V^*(s')]$$

Run a fixed point iteration to find V*:

- Start with initial guess V^0 .
- For i = 1, ..., z:
 - For $s \in \mathcal{S}$:
 - $V^{i}(s) = \max_{a} \sum_{s',r} \cdot \Pr(s',r \mid s,a)[r + \gamma V^{i-1}(s')]$

Can be shown to converge in roughly $z = \frac{1}{1-\gamma}$ iterations. What is the computational cost of each iteration?

Full information: We know S, A, the transition function P and reward function R.

Reinforcement Learning setting: We do not know *P* or *R*, but we can repeatedly play the MDP, running whatever policy we like.

- Model-based RL methods essentially try to learn P and R very accurately and then find Π^* via a method like value iteration. Require a lot of samples of the MDP.
- Model-free RL methods try to learn Π* without necessarily obtaining an accurate model of the world – i.e. without explicitly learning P and R.

Q FUNCTION

Another important function:

• Q-function: $Q^{\Pi}(s, a) = \mathbb{E}_{\Pi, s_0=s, a_0=a} \sum_{t\geq 0} \gamma^t r_t$. Measures the expected return if we start in state *s*, play action *a*, and then follow policy Π .

$$Q^*(s, a) = \max_{\Pi} Q^{\Pi}(s, a) = Q^{\Pi^*}(s, a).$$



Q FUNCTION

$$Q^*(s,a) = \max_{\Pi} \mathbb{E}_{\Pi, s_0=s, a_0=a} \sum_{t\geq 0} \gamma^t r_t.$$

If we knew the function Q^* , we would immediately know an optimal policy. Whenever we're in state *s*, we should always play action $a^* = \arg \max_a Q^*(s, a)$.



Q has more parameters than *V*, but you can use it to determine an optimal policy without knowing transition probabilities.

Q* also satisfies a Bellman equation:

$$Q^*(s,a) = \mathbb{E}[R(s,a)] + \gamma \mathbb{E}_{s' \sim P(s,a)} \max_{a'} Q^*(s',a').$$

Bellman equation:

$$Q^*(s,a) = \mathbb{E}[R(s,a)] + \gamma \mathbb{E}_{s' \sim P(s,a)} \max_{a'} Q^*(s',a').$$

Again use fixed point iteration to find Q^* . Let Q^{i-1} be our current guess for Q^* and suppose we are at some state *s*, *a*.

$$Q^{i}(s,a) = \mathbb{E}[R(s,a)] + \gamma \mathbb{E}_{s' \sim P(s,a)} \max_{a'} Q^{i-1}(s',a')$$

In reality, drop expectations and use a learning rate lpha

$$Q^{i}(s,a) = (1-\alpha)Q^{i}(s,a) + \alpha \left(R(s,a) + \gamma \max_{a'} Q^{i-1}(s',a')\right)$$

Q LEARNING

How do we choose states s and a to make the update for? In principal you can do anything you want! E.g. choose some policy Π and run:

- Initialize Q⁰ (e.g. all zeros)
- Start at s, play action $a = \Pi(s)$, observe reward R(s, a).
- For i = 1, ..., z
 - $Q^{i}(s,a) = (1-\alpha)Q^{i}(s,a) + \alpha \left(R(s,a) + \gamma \max_{a'} Q^{i-1}(s',a')\right)$
 - \cdot s \leftarrow P(s, a)
 - · $a \leftarrow \Pi(s)$

(restart if we reach a terminating state)

Q-learning is considered an **off-policy** RL method because it runs a policy Π that is not necessarily related to its current guess for an optimal policy, which in this case would be $\Pi(s) = \max_a Q^i(s, a)$ at time *i*. For small enough α , Q-learning converges to Q^* as long as we follow a policy Π that visits every start (s, a) with non-zero probability.

Mild condition, but exact choice of Π matters for convergence rate.

- Random: At state s, choose a random action a.
- **Greedy:** At state s, choose $\arg \max_a Q^i(s, a)$. I.e. the current guess for the best action.

	end +1
	end -1
start	

Random can be wasteful. Spend time improving parts of *Q* that aren't relevant to optimal play. **Greedy** can cause you to zero in on a locally optimal policy without learning new strategies.

Possible choices for Π :

- Random: At state s, choose a random action a.
- **Greedy:** At state *s*, choose $\arg \max_a Q^i(s, a)$. I.e. the current guess for the best action.
- ϵ -Greedy: At state s, choose $\arg \max_a Q^i(s, a)$ with probability 1ϵ and a random action with probability ϵ .



Exploration-exploitation tradeoff. Increasing ϵ = more **exploration**.

Another issue: Even writing down Q^* is intractable... This is a function over |S||A| possible inputs. Even for relatively simple games, |S| is gigantic...

Back of the envelope calculations:

- Tic-tac-toe: $3^{(3\times3)} \approx 20,000$
- Chess: $\approx 10^{43} < 28^{64}$ (due to Claude Shannon).
- **Go:** $3^{(19 \times 19)} \approx 10^{171}$.
- Atari: $128^{(210 \times 160)} \approx 10^{71,000}$.

Number of atoms in the universe: $\approx 10^{82}$.

Learn a **simpler** function $Q(s, a, \theta) \approx Q^*(s, a)$ parameterized by a small number of parameters θ .

Example: Suppose our state can be represented by a vector in \mathbb{R}^d and our action *a* by an integer in $1, \ldots, |\mathcal{A}|$. We could use a linear function where θ is a small matrix:



Learn a **simpler** function $Q(s, a, \theta) \approx Q^*(s, a)$ parameterized by a small number of parameters θ .

Example: Could also use a (deep) neural network.



DeepMind: "Human-level control through deep reinforcement learning", Nature 2015.

If $Q(s, a, \theta)$ is a good approximation to $Q^*(s, a)$ then we have an approximately optimal policy: $\tilde{\Pi}^*(s) = \arg \max_a Q(s, a, \theta)$.

- Start in state s₀.
- For t = 1, 2, ...
 - $a^* = \arg \max_a Q(s, a, \theta)$
 - $s_t \sim P(s_{t-1}, a^*)$

How do we find an optimal θ ? If we knew $Q^*(s, a)$ could use supervised learning, but the true Q function is infeasible to compute.

Q-LEARNING W/ FUNCTION APPROXIMATION

Find θ which satisfies the Bellman equation:

$$Q^*(s, a) = \mathbb{E}_{s' \sim P(s, a)} \left[R(s, a) + \gamma \max_{a'} Q^*(s', a') \right]$$
$$Q(s, a, \theta) \approx \mathbb{E}_{s' \sim P(s, a)} \left[R(s, a) + \gamma \max_{a'} Q(s, a, \theta) \right].$$

Should be true for all a, s. Should also be true for $a, s \sim D$ for any distribution D:

$$\mathbb{E}_{s,a\sim\mathcal{D}}Q(s,a,\theta)\approx\mathbb{E}_{s,a\sim\mathcal{D}}\mathbb{E}_{s'\sim\mathcal{P}(s,a)}\left[R(s,a)+\gamma\max_{a'}Q(s,a,\theta)\right].$$

Loss function:

$$L(\theta) = \mathbb{E}_{s,a \sim D} (y - Q(s, a, \theta))^2$$

where $y = \mathbb{E}_{s' \sim P(s,a)} [R(s,a) + \gamma \max_{a'} Q(s',a',\theta)].$

Q-LEARNING W/ FUNCTION APPROXIMATION

Minimize loss with gradient descent:

$$\nabla L(\theta) = \mathbb{E}_{s,a \sim \mathcal{D}} \left[-2\nabla Q(s,a,\theta) \cdot \left[y - Q(s,a,\theta) \right] \right]$$

In practice use stochastic gradient:

$$\nabla L(\theta, s, a) = -2 \cdot \nabla Q(s, a, \theta) \cdot \left[R(s, a) + \gamma \max_{a'} Q(s', a', \theta) - Q(s, a, \theta) \right]$$

- Initialize θ_0
- For *i* = 0, 1, 2, ...
 - Run policy Π to obtain s, a and s' ~ P(s, a)
 - Set $\theta_{i+1} = \theta_i \eta \cdot \nabla L(\theta_i, s, a)$

 η is a learning rate parameter.
Again, the choice of Π matters a lot. **Random play** can be wastefully, putting effort into approximating Q^* well in parts of the state-action space that don't actually matter for optimal play. ϵ -greedy approach is much more common:

• Initialize s₀.

For
$$t = 0, 1, 2, ...,$$

 $\cdot a_i = \begin{cases} \arg \max_a Q(s_t, a, \theta_{curr}) & \text{with probability } (1 - \epsilon) \\ \text{random action} & \text{with probability } \epsilon \end{cases}$

Lots of other details we don't have time for! References:

- Original DeepMind Atari paper: https://www.cs.toronto.edu/~vmnih/docs/dqn.pdf, which is very readable.
- Stanford lecture video: https://www.youtube.com/watch?v=lvoHnicueoE and slides: http://cs231n.stanford.edu/slides/2017/ cs231n_2017_lecture14.pdf

Important concept we did not cover: experience replay.



https://www.youtube.com/watch?v=V1eYniJ0Rnk