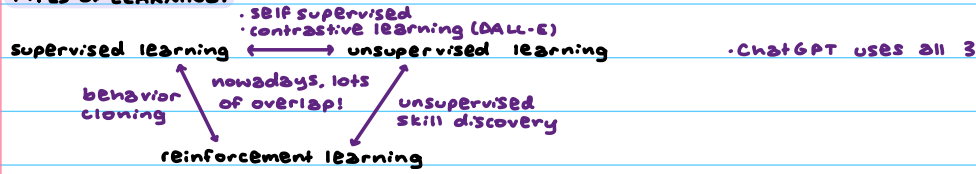


6.390 NOTES - SHEN SHEN

TYPES OF LEARNING:



SUPERVISED LEARNING:

REGRESSION: predicting a continuous number

ex: predicting city's energy usage

TRAINING DATA:  $D_{train} := \{(x^1, y^1), \dots, (x^d, y^d)\}$

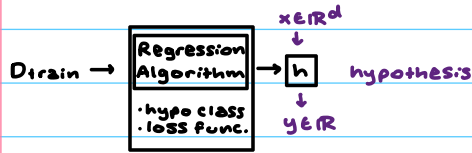
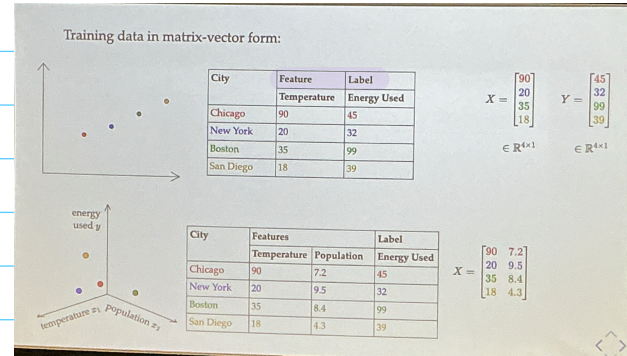
feature vector ( $x$ )      label ( $y$ )

$x^i = \begin{bmatrix} x_1^i \\ x_2^i \\ \vdots \\ x_d^i \end{bmatrix} \in \mathbb{R}^d$        $y^i \in \mathbb{R} \leftarrow$  always scalar

- input variable (what model knows)
- output / target var (what model is trying to predict)

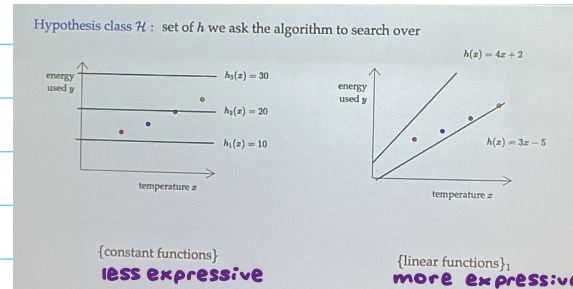
$n$  data pts, each w/  $d$  dimensional features & scalar label

$n \times d$   $x$ ,  $n \times 1$   $y$



HYPOTHESIS: function that can take in new  $x$  & output new  $y$

- want to make good predictions
- HYPOTHESIS CLASS  $\mathcal{H}$ : set of  $h$  we ask algorithm to search over

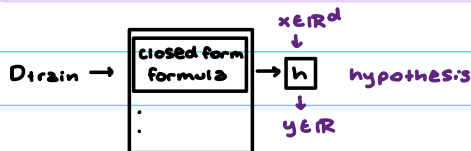


LOSS:  $L(h(x^i), y^i)$  how good hypothesis is compared to actual outcome

TRAINING ERROR:  $\mathcal{E}_{train}(h) = \frac{1}{n} \sum_{i=1}^n L(h(x^i), y^i)$  ← average loss on training data

TEST ERROR:  $\mathcal{E}_{test}(h) = \frac{1}{n'} \sum_{i=n+1}^{n+n'} L(h(x^i), y^i)$  ← average loss on  $n'$  unseen test data pts

ORDINARY LEAST SQUARES REGRESSION:



LINEAR HYPOTHESIS CLASS:  $h(x, \theta) = [\theta_0, \theta_1, \dots, \theta_d] \cdot \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_d \end{bmatrix} = \theta^T x$

params

Features

SQUARED LOSS FUNC:  $L(h(x^i), y^i) = (\theta^T x^i - y^i)^2$

DERIVE OLS:

- write training error  $J(\theta)$  in scalar form
- rearrange into matrix-vector form  $J(\theta) = \frac{1}{2} (X\theta - Y)^T (X\theta - Y)$  ← objective
- set gradient  $\nabla_{\theta} J$  to zero
- solve optimal params  $\theta^* = (X^T X)^{-1} X^T Y$  ← closed form soln: when well defined,  $\theta^*$  is unique minimizer of  $J(\theta)$

1. Write out training error:

e.g.:

| City      | Features  |         | Label        |
|-----------|-----------|---------|--------------|
|           | Temp (°F) | Pop (M) | Energy (kWh) |
| Chicago   | 90        | 7.2     | 45           |
| New York  | 20        | 9.5     | 32           |
| Boston    | 35        | 8.4     | 99           |
| San Diego | 18        | 4.3     | 39           |

avg. training data

$$J(\theta) = \frac{1}{2} [(\theta_1 + 90 + \theta_2 - 45)^2 + (\theta_1 + 20 + \theta_2 - 9.5 - 32)^2 + (\theta_1 + 35 + \theta_2 - 8.4 - 99)^2 + (\theta_1 + 18 + \theta_2 - 4.3 - 39)^2]$$

$$\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix} \in \mathbb{R}^{2 \times 1}$$

• Q: What kind of function is  $J(\theta)$ ?  
quadratic (2nd order)

• Q: what  $J(\theta)$  look like?  
typically a bowl



2. Rearrange training error into matrix-vector form

| City      | Features  |         | Label        |
|-----------|-----------|---------|--------------|
|           | Temp (°F) | Pop (M) | Energy (kWh) |
| Chicago   | 90        | 7.2     | 45           |
| New York  | 20        | 9.5     | 32           |
| Boston    | 35        | 8.4     | 99           |
| San Diego | 18        | 4.3     | 39           |

$$X = \begin{bmatrix} 90 & 7.2 \\ 20 & 9.5 \\ 35 & 8.4 \\ 18 & 4.3 \end{bmatrix} \quad Y = \begin{bmatrix} 45 \\ 32 \\ 99 \\ 39 \end{bmatrix} \quad \theta = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix}$$

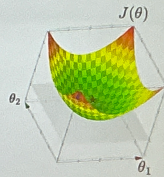
$$\in \mathbb{R}^{4 \times 2} \quad \in \mathbb{R}^{4 \times 1} \quad \in \mathbb{R}^{2 \times 1}$$

$$J(\theta) = \frac{1}{2} [(\theta_1 + 90 + \theta_2 - 45)^2 + (\theta_1 + 20 + \theta_2 - 9.5 - 32)^2 + (\theta_1 + 35 + \theta_2 - 8.4 - 99)^2 + (\theta_1 + 18 + \theta_2 - 4.3 - 39)^2]$$

$$J(\theta) = \frac{1}{2} (X\theta - Y)^T (X\theta - Y)$$

3. Get the gradient  $\nabla_{\theta} J \stackrel{!}{=} 0$

$$\nabla_{\theta} J = \begin{bmatrix} \partial J / \partial \theta_1 \\ \vdots \\ \partial J / \partial \theta_d \end{bmatrix} = \frac{\partial}{\partial \theta} (X^T X \theta - X^T Y)$$



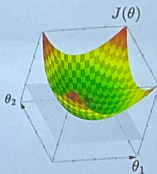
4. Set the gradient  $\nabla_{\theta} J \stackrel{!}{=} 0$

$$\Rightarrow \theta^* = (X^T X)^{-1} X^T Y$$

- Typically,  $J(\theta)$  "curves up"
- The minimizer of  $J(\theta)$  necessarily has a gradient zero.

The beauty of

$$\theta^* = (X^T X)^{-1} X^T Y$$



- When  $\theta^*$  is well defined, it's the unique minimizer of  $J(\theta)$
- Closed-form solution, does not feel like "training"
- Very rare case where we get a general and clean solution with nice theoretical guarantee.

## 2/9 REGULARIZATION & CROSS-VALIDATION:

Recall

Let

$$X = \begin{bmatrix} -x^{(1)\top} \\ \vdots \\ -x^{(n)\top} \end{bmatrix} = \begin{bmatrix} x_1^{(1)} & \dots & x_d^{(1)} \\ \vdots & \ddots & \vdots \\ x_1^{(n)} & \dots & x_d^{(n)} \end{bmatrix} \quad Y = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix} \quad \theta = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_d \end{bmatrix}$$

$\in \mathbb{R}^{n \times d} \quad \in \mathbb{R}^{n \times 1} \quad \in \mathbb{R}^{d \times 1}$

Then

$$J(\theta) = \frac{1}{2} (X\theta - Y)^T (X\theta - Y) \in \mathbb{R}^{1 \times 1}$$

$$\theta^* = (X^T X)^{-1} X^T Y \in \mathbb{R}^{d \times 1} \quad \leftarrow \text{only works if } X \text{ \& } Y \text{ are vectors}$$

$$\theta^* = [2.00, 3.00]^T \quad J(\theta^*) = 0.00$$

Feature (d=2) and Label

$$\begin{bmatrix} 1 \\ 0 \end{bmatrix}, 2 \quad \begin{bmatrix} 0 \\ 1 \end{bmatrix}, 3 \quad \begin{bmatrix} 1 \\ 1 \end{bmatrix}, 5$$

hypothesis  $h = 2x_1 + 3x_2$

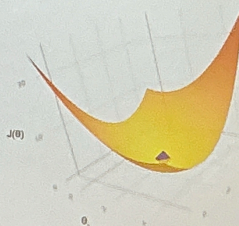
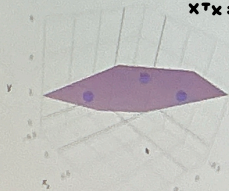
Training error  $J(\theta_1, \theta_2)$

$$J(\theta) = \frac{1}{2} [(\theta_1 - 2)^2 + (\theta_2 - 3)^2 + (\theta_1 + \theta_2 - 5)^2]$$

$$X = \begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{bmatrix} \quad Y = \begin{bmatrix} 2 \\ 3 \\ 5 \end{bmatrix}$$

$$X^T X = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} \quad X^T Y = \begin{bmatrix} 3 \\ 8 \end{bmatrix}$$

$$\theta^* = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix}^{-1} \begin{bmatrix} 3 \\ 8 \end{bmatrix} = \begin{bmatrix} 2 \\ 3 \end{bmatrix}$$



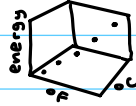
PROBLEMS: closed-form formula UNDEFINED

a)  $n > d$  (less data than feature we're trying to fit):



-ex: genomics, NLP

b) linearly-dependent features



-ex: temp °F/°C, age/birth-yr

closed form formula:  $\theta^* = (X^T X)^{-1} X^T Y$  undefined

optimal soln:  $\infty$  many optimal  $\theta^*$   $\rightarrow$  not enough info to pin down a unique soln.

mathematically,

$(X^T X)$  is singular

$\Updownarrow$

$X$  is not full column rank

MM 2

$Ax$  and  $Ay$  are linear combinations of columns of  $A$ .

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \end{bmatrix} \begin{bmatrix} x_1 & y_1 \\ x_2 & y_2 \end{bmatrix} = A \begin{bmatrix} x & y \end{bmatrix} = \begin{bmatrix} Ax & Ay \end{bmatrix}$$

demo (a)  $n < d$ : 1 sample, 2 features

$$X^T X = \begin{bmatrix} 2 & 3 \\ 3 & 2 \end{bmatrix} = \begin{bmatrix} 4 & 6 \\ 6 & 9 \end{bmatrix}$$

$\begin{bmatrix} 3 \\ 2 \end{bmatrix} \begin{bmatrix} 2 & 3 \end{bmatrix}$

demo (b) Collinear:  $x_2 = 1.5 \cdot x_1$

$$X^T X = \begin{bmatrix} 2 & 4 & 6 \\ 3 & 6 & 9 \end{bmatrix} \begin{bmatrix} 2 & 3 \\ 4 & 6 \\ 6 & 9 \end{bmatrix} = \begin{bmatrix} 56 & 84 \\ 84 & 126 \end{bmatrix}$$

$J(\theta)$

$J(\theta)$

When  $X$  is not full column rank

- $J(\theta)$  has a "flat" bottom
- This  $\theta^*$  formula is not well-defined
- Infinitely many optimal hyperplanes

Typically,  $X$  is full column rank

- $J(\theta)$  "curves up" everywhere
- $\theta^* = (X^T X)^{-1} X^T Y$
- unique optimal hyperplane

formula itself isn't wrong, data is trouble-making  $X^T X$  becoming more invertible

$\min \lambda$  of  $X^T X > 0$  (we want all positive semidefinites)

sits in middle (very sensitive info)

ex: assume  $n=1$  &  $y=1$ , then  $\theta^* = \frac{1}{x}$

if data is  $(x, 1) = (0.002, 1) \rightarrow \theta^* = 500$



if data is  $(x, 1) = (-0.002, 1) \rightarrow \theta^* = -5000$

when  $X^T X$  is almost singular

technically,  $\theta^* = (X^T X)^{-1} X^T Y$  exists and gives the unique optimal hyperplane

practically,  $\theta^*$  tends to have huge magnitude

$\theta^*$  tends to be very sensitive to the small changes in the data

lots of other  $\theta$ s fit the training data almost equally well

RIDGE REGRESSION:  $J_{\text{ridge}}(\theta) = \frac{1}{n} (X\theta - y)^T (X\theta - y) + \lambda \|\theta\|^2$

MSE on training data
penalty on param magnitude

$\lambda > 0$  controls how heavily we penalize magnitude relative to MSE

- many  $\theta$  give similar loss, but some have huge magnitude  $\rightarrow$  unstable!
- small change in  $x \rightarrow$  wildly different prediction

IDEA: penalize large  $\theta$  in our obj (explicit regularization)

$$\theta_{\text{ridge}}^* = (X^T X + n\lambda I)^{-1} X^T Y$$

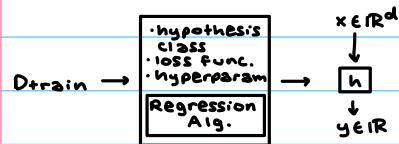
for  $\lambda > 0$ :  $X^T X + n\lambda I$  is always invertible, so  $\theta_{\text{ridge}}^*$  always exists and is unique

ex:  $\lambda = 0 \rightarrow$  no penalty - reduces to OLS

$\lambda = 1000 \rightarrow$  huge penalty - forces  $\theta \approx 0$

$\lambda = -100 \rightarrow$  rewards large  $\theta$  - counter-productive  $\leftarrow \lambda > 0$  required

$\lambda$  is hyperparameter:



- affects learning outcome, not learned by algo.

we need to choose hyperparams like  $\lambda$ :

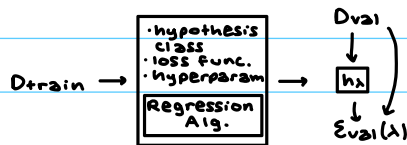
- can't use training error (would always pick  $\lambda = 0$ )
- can't use test error (we don't have test data)

hold-out some training data  $\{(x^1, y^1), \dots, (x^n, y^n)\}$

use  $D_{\text{val}}$  to evaluate how

|                    |                  |
|--------------------|------------------|
| $D_{\text{train}}$ | $D_{\text{val}}$ |
|--------------------|------------------|

good  $\lambda$  hyperparam is



for each  $\lambda \in \{0.1, 1, 10\}$ :

- train on  $D_{\text{train}}$  with  $\lambda$
- compute  $E_{\text{val}}(\lambda)$  on  $D_{\text{val}}$

return  $\arg \min_{\lambda} E_{\text{val}}(\lambda) \rightarrow$  in this example, compare  $E_{\text{val}}(0.1), (1), (10)$ , return smallest validation error

### Cross-validation

for each  $\lambda \in \{0.1, 1, 10\}$ :

for  $i = 1, \dots, 5$ :

train  $h_i$  on  $\mathcal{D} \setminus \mathcal{D}_i$  with  $\lambda$

$\mathcal{E}_i$  = error on  $\mathcal{D}_i$

$$E_{\text{val}}(\lambda) = (\mathcal{E}_1 + \dots + \mathcal{E}_5) / 5$$

return  $\lambda^* = \arg \min_{\lambda} E_{\text{val}}(\lambda)$

$\lambda \in \{0.1, 1, 10\}$ :



$$E_{\text{val}}(\lambda) = (\mathcal{E}_1 + \mathcal{E}_2 + \mathcal{E}_3 + \mathcal{E}_4 + \mathcal{E}_5) / 5$$

- divide data
- train chunks
- choose  $\lambda$
- find  $\lambda^*$

How many hypotheses trained in this example to pick  $\lambda^*$ ?

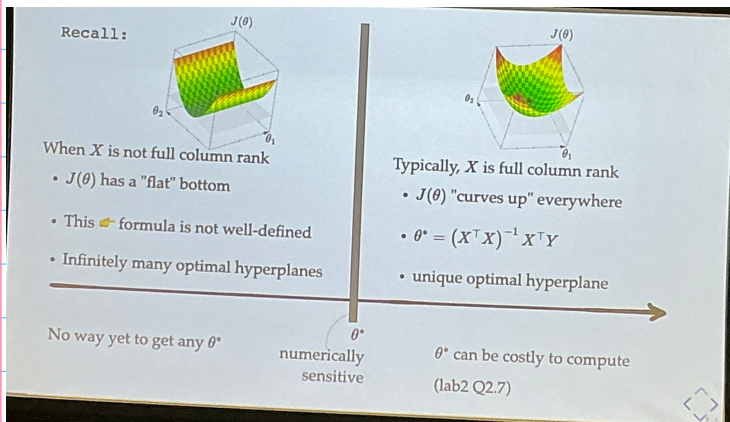
$$\theta_{\text{final}}^* = (X^T X + n\lambda^* I)^{-1} X^T Y$$

data from all of  $\mathcal{D}_1$  to  $\mathcal{D}_5$

# Summary

- When  $X^T X$  is singular or ill-conditioned, OLS is undefined or overfits.
- Regularization combats overfitting by penalizing large  $\theta$ .
- Ridge regression adds  $\lambda \|\theta\|^2$  to the objective — still has a closed-form solution.
- $\lambda$  is a hyperparameter that trades off fit vs. regularization.
- Validation and cross-validation provide principled ways to choose  $\lambda$ .

## 2/17 GRADIENT DESCENT METHODS (more general & efficient algo to train our ML model)



### GRADIENT DESCENT:

GRADIENT OF:  $\mathbb{R}^m \rightarrow \mathbb{R}^m$  is defined @ point  $p = (x_1, \dots, x_m)$  as:  $\nabla f(p) = \begin{bmatrix} \frac{\partial f}{\partial x_1}(p) \\ \vdots \\ \frac{\partial f}{\partial x_n}(p) \end{bmatrix}$

1: gradient generalizes concept of derivative to multiple dims.

2: gradient's dimensionality always matches function input

3: can be symbolic or numeric. evaluating symbolic gives numeric

4: gradient points in direction of (steepest) increase in function value

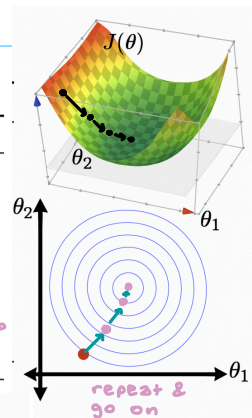
5: gradient @ function minimizer is necessarily zero

### GD ALGORITHM:

hyperparameters      initial guess of parameters      precision

Algorithm 1 Gradient Descent ( $\theta_{\text{init}}, \eta, J, \epsilon$ )

- 1: Initialize  $\theta^{(0)} = \theta_{\text{init}}$  • learning rate
- 2: Initialize  $t = 0$
- 3: repeat
- 4:  $t = t + 1$  • scaling factor
- 5:  $\theta^{(t)} = \theta^{(t-1)} - \eta \nabla_{\theta} J(\theta^{(t-1)})$
- 6: until  $|J(\theta^{(t)}) - J(\theta^{(t-1)})| < \epsilon$  ← precision gap
- 7: return  $\theta^{(t)}$



Algorithm 1 Gradient Descent ( $\theta_{\text{init}}, \eta, J, \epsilon$ )

- 1: Initialize  $\theta^{(0)} = \theta_{\text{init}}$
- 2: Initialize  $t = 0$
- 3: repeat
- 4:  $t = t + 1$
- 5:  $\theta^{(t)} = \theta^{(t-1)} - \eta \nabla_{\theta} J(\theta^{(t-1)})$
- 6: until  $|J(\theta^{(t)}) - J(\theta^{(t-1)})| < \epsilon$  ← objective improvement is nearly zero.
- 7: return  $\theta^{(t)}$

Other possible stopping criterion for line 6:

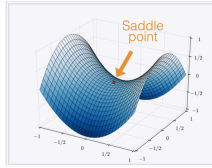
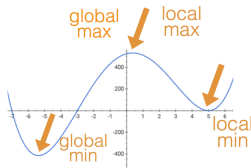
- Small parameter change:  $\|\theta^{(t)} - \theta^{(t-1)}\| < \epsilon$ , or
- Small gradient norm:  $\|\nabla_{\theta} J(\theta^{(t-1)})\| < \epsilon$

#tracking improvement over time

**GRADIENT DESCENT PROPERTIES:**

When minimizing a function, we aim for a global minimizer.

At a global minimizer  $\Rightarrow$  the gradient vector is zero  $\Leftarrow$  gradient descent can achieve this (to arbitrary precision)



at a global minimizer  $\Leftarrow$   $\begin{cases} \text{the gradient vector is zero} \\ \text{objective function is convex} \end{cases}$   $\Leftarrow$  gradient desc. can achieve this to arbitrary prec.

**CONVEX:** any line segment connecting graph  $f$  lies above or on the graph

• optimization theory guarantees convex funcs. converge / efficiency

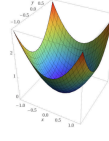
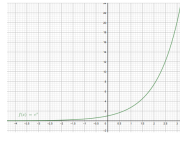
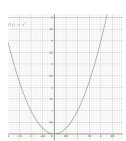


• Jmse always convex

• Jridge with  $\lambda > 0$  always strongly convex

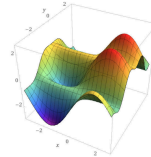
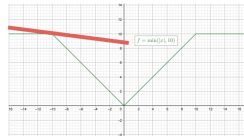
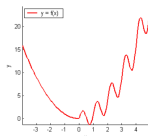
More examples

Convex functions



any 2 pts that are connected on graph, the line segment lives above the curve.

Non-convex functions

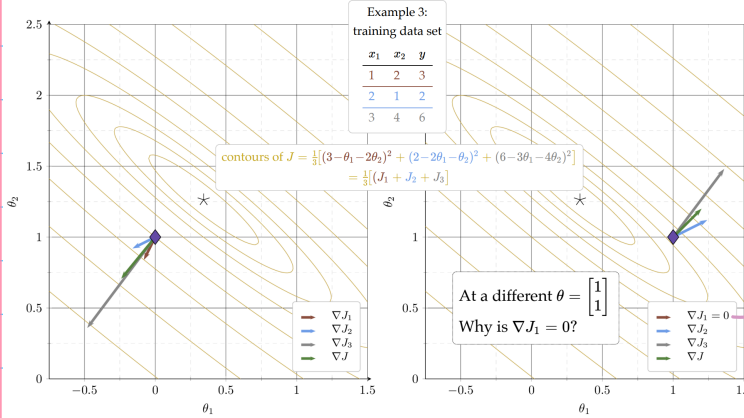


**GRADIENT DESCENT PERFORMANCE:**

Gradient Descent Performance

- Assumptions:
  - $f$  is sufficiently "smooth"  $\rightarrow$  if violated, might not have gradient & can't run gradient descent
  - $f$  is convex  $\rightarrow$  if violated, might get stuck @ saddle point
  - $f$  has at least one global minimum  $\rightarrow$  if violated, may not terminate/converge (ex: a line) ✗
  - Run gradient descent for sufficient iterations
  - $\eta$  is sufficiently small
- Conclusion:
  - Gradient descent converges arbitrarily close to a global minimizer of  $f$ .

**STOCHASTIC GRADIENT DESCENT (SGD):**



green arrow is the 1/3 avg.

pt (1,1) is best fit param already, so grad = 0 (no room for improvement)

**Gradient of an ML objective**

In general,

- the MSE of a linear hypothesis:  $J(\theta) = \frac{1}{n} \sum_{i=1}^n (\theta^T x^{(i)} - y^{(i)})^2$
- and its gradient w.r.t.  $\theta$ :  $\nabla_{\theta} J(\theta) = \frac{1}{n} \sum_{i=1}^n 2(\theta^T x^{(i)} - y^{(i)}) x^{(i)}$

An ML objective function is a finite sum

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n J_i(\theta)$$

and its gradient w.r.t.  $\theta$ :

$$\nabla_{\theta} J(\theta) = \nabla_{\theta} \left( \frac{1}{n} \sum_{i=1}^n J_i(\theta) \right) = \frac{1}{n} \sum_{i=1}^n \nabla_{\theta} J_i(\theta)$$

☺ (gradient of the sum) = (sum of the gradient)

**Gradient of an ML objective**

In general,

- An ML objective function is a finite sum

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n J_i(\theta)$$

loss incurred on a single  $i^{\text{th}}$  data point

- and its gradient w.r.t.  $\theta$ :  $\nabla_{\theta} J(\theta) = \nabla_{\theta} \left( \frac{1}{n} \sum_{i=1}^n J_i(\theta) \right) = \frac{1}{n} \sum_{i=1}^n \nabla_{\theta} J_i(\theta)$

need to add  $n$  of these, each  $\nabla_{\theta} J_i \in \mathbb{R}^d$   
 = one for each pt

Costly in practice! gradient info from the  $i^{\text{th}}$  data point's loss

SGD is rougher but less costly

**Algorithm 1 Gradient Descent**( $\theta_{\text{init}}, \eta, J, \epsilon$ )

- 1: Initialize  $\theta^{(0)} = \theta_{\text{init}}$
- 2: Initialize  $t = 0$
- 3: **repeat**
- 4:  $t = t + 1$
- 5:  $\theta^{(t)} = \theta^{(t-1)} - \eta \nabla_{\theta} J(\theta^{(t-1)})$
- 6: **until**  $|J(\theta^{(t)}) - J(\theta^{(t-1)})| < \epsilon$
- 7: **return**  $\theta^{(t)}$

**Algorithm 2 Stochastic Gradient Descent**( $\theta_{\text{init}}, \eta, \{J_i\}_{i=1}^n, \epsilon$ )

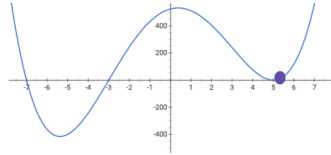
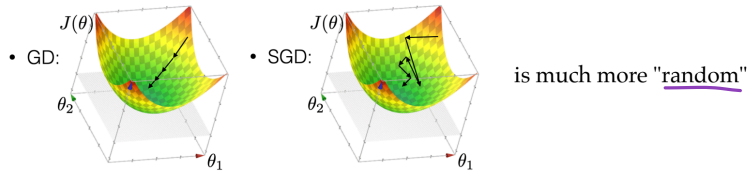
- 1: Initialize  $\theta^{(0)} = \theta_{\text{init}}$
- 2: Initialize  $t = 0$
- 3: **repeat**
- 4:  $t = t + 1$
- 5:  $i = \text{randomly sample from } \{1, \dots, n\}$
- 6:  $\theta^{(t)} = \theta^{(t-1)} - \eta \nabla_{\theta} J_i(\theta^{(t-1)})$
- 7: **until**  $|J(\theta^{(t)}) - J(\theta^{(t-1)})| < \epsilon$
- 8: **return**  $\theta^{(t)}$

$$\nabla_{\theta} J(\theta) = \frac{1}{n} \sum_{i=1}^n \nabla_{\theta} J_i(\theta) \approx \nabla_{\theta} J_i(\theta)$$

← randomly chosen sample is the true gradient

Compared with GD, SGD:

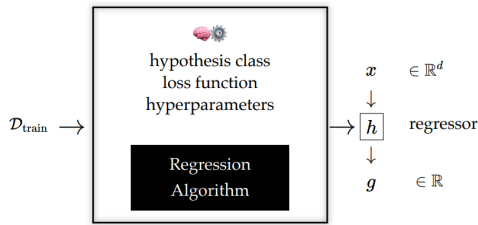
$$\nabla_{\theta} J(\theta) = \frac{1}{n} \sum_{i=1}^n \nabla_{\theta} J_i(\theta) \approx \nabla_{\theta} J_i(\theta) \quad \text{is more efficient}$$



2/23 **LINEAR CLASSIFICATION**

- binary: separator, normal vector
- logistic: sigmoid, neg. log-likelihood loss
- multi-class: softmax, one-hot encoding, cross-entropy loss

Recap:



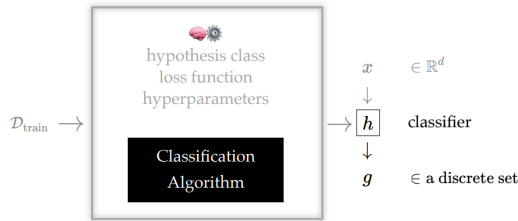
"Learn" a model →

train, optimize, tune, adapt ...  
adjusting/updating/finding  $\theta$   
gradient based

↓ "Use" a model

predict, test, evaluate, infer ...  
plug in the  $\theta$  found  
no gradients involved

Today:



\*intuitively, seems easier since output is discrete set, but it's actually more difficult

{0, 1}

{😊 😞}

["good", "better", "best", ...]

["Fish", "Grizzly", "Chameleon", ...]

**LINEAR (binary) CLASSIFIERS:**

|                          | linear regressor  | linear binary classifier   |
|--------------------------|---|--|
| features                 | $x \in \mathbb{R}^d$  | $x \in \mathbb{R}^d$   |
|                          | $y \in \mathbb{R}$  | $y \in \{0, 1\}$   |
| parameters               | $\theta \in \mathbb{R}^d, \theta_0 \in \mathbb{R}$  |  |
| linear combo             | $\theta^T x + \theta_0 = z \leftarrow \text{linear combo.}$                                     |  |
| predict                  | $g = z$   | $g = \begin{cases} 1 & \text{if } z > 0 \\ 0 & \text{otherwise} \end{cases}$ $g$ is "flat" and discrete in $\theta$                                |
| loss $\mathcal{L}(g, y)$ | $\mathcal{L}_{\text{squared}} = (g - y)^2$  | $\mathcal{L}_{01} = \begin{cases} 0 & \text{if } g = y \\ 1 & \text{otherwise} \end{cases}$ $\mathcal{L}_{01}(g, y)$ is "flat" and discrete in $g$ |
| optimize method          | <ul style="list-style-type: none"> <li>closed-form formula</li> <li>gradient descent</li> </ul> | training error almost "flat" w.r.t $\theta$ , gradient gives very little info  |

← 2 options

separator: @ 0 threshold  
@ always points @ + feature space

natural loss func. choice:  $g = \text{step}(z) = \text{step}(\theta^T x + \theta_0)$

$\mathcal{L}_{01}(g, y) = \begin{cases} 0 & \text{if } g = y \\ 1 & \text{o.w.} \end{cases}$  → VERY hard to optimize. flat almost everywhere (zero grad.), "jumps" elsewhere

**LINEAR LOGISTIC (binary) CLASSIFIERS:**

**Sigmoid  $\sigma(\cdot)$ :** confidence or estimated likelihood that  $x$  belongs to the positive class

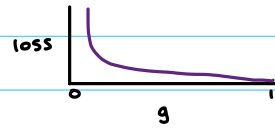
- monotonic, elegant gradient 'smoothed version'

-  $\theta, \theta_0$  can flip, squeeze, expand, or shift the  $\sigma(x)$  graph horizontally

|              | linear binary classifier   | linear logistic binary classifier  |
|--------------|--|--|
| features     |  | $x \in \mathbb{R}^d$   |
| parameters   |  | $\theta \in \mathbb{R}^d, \theta_0 \in \mathbb{R}$                                 |
| linear combo |  | $\theta^T x + \theta_0 = z$  |
| predict      |  |  |
|              | $\begin{cases} 1 & \text{if } z > 0 \\ 0 & \text{otherwise} \end{cases}$ | $\begin{cases} 1 & \text{if } \sigma(z) > 0.5 \\ 0 & \text{otherwise} \end{cases}$ |

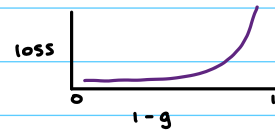
Sigmoid  $\sigma(\cdot)$ : confidence or estimated likelihood that  $x$  belongs to the positive class

loss func:  $\mathcal{L}_{\text{ML}}(g, y) = -\log g$  <sup>likelihood</sup>



- Smooth
- When  $g$  almost 1  $\rightarrow$  little loss
- $g$  close to 0  $\rightarrow$  big loss

$1-g = 1-\sigma(\cdot)$ : model's predicted likelihood that  $x$  belongs to the negative class



- $\log(1-g)$
- low loss = low

**linear logistic binary classifier**

features:  $x \in \mathbb{R}^d$   
parameters:  $\theta \in \mathbb{R}^d, \theta_0 \in \mathbb{R}$

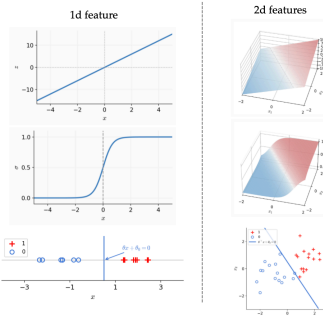
the logit  $z$ :  $z = \theta^T x + \theta_0$

apply sigmoid:

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

Predict 1 if  $\sigma(z) > 0.5$ , else 0.

$$\sigma(z) = 0.5 \iff z = 0 \iff \theta^T x + \theta_0 = 0$$



separator is linear in feature  $x$ !

$$\mathcal{L}_{\text{ML}}(g, y) = \begin{cases} -\log(g) & \text{if } y = 1 \\ -\log(1-g) & \text{if } y = 0 \end{cases} \iff -[y \log g + (1-y) \log(1-g)]$$

$\leftarrow$  encourage  $g$  closer to 1  
 $\leftarrow$   $g$  closer to 0

• When  $y = 1$ :  $-[y \log g + (1-y) \log(1-g)] = -\log g$

• When  $y = 0$ :  $-[y \log g + (1-y) \log(1-g)] = -\log(1-g)$

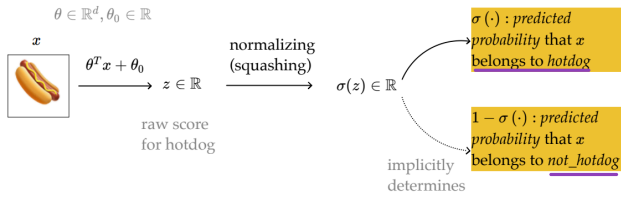
Read as:  $\sum$  (true label for class  $k$ )  $\cdot -\log$ (predicted prob of class  $k$ ).

Since  $y \in \{0, 1\}$ , only the true class's term survives.

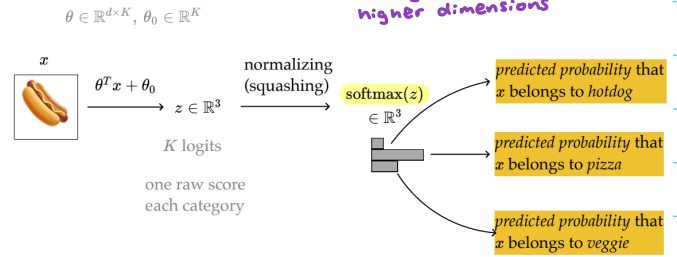
|              | linear binary classifier   | linear logistic binary classifier  |
|--------------|--|--|
| features     |  | $x \in \mathbb{R}^d$   |
| label        |  | $y \in \{0, 1\}$   |
| parameters   |  | $\theta \in \mathbb{R}^d, \theta_0 \in \mathbb{R}$                                       |
| linear combo |  | $\theta^T x + \theta_0 = z$  |
| predict      | $\begin{cases} 1 & \text{if } z > 0 \\ 0 & \text{otherwise} \end{cases}$ | $\begin{cases} 1 & \text{if } g = \sigma(z) > 0.5 \\ 0 & \text{otherwise} \end{cases}$   |
| loss         | $\begin{cases} 0 & \text{if } g = y \\ 1 & \text{otherwise} \end{cases}$ | $\begin{cases} -\log(g) & \text{if } y = 1 \\ -\log(1-g) & \text{if } y = 0 \end{cases}$ |
|              |  | $\iff -[y \log g + (1-y) \log(1-g)]$   |
| optimize via | NP-hard to learn   | gradient descent   |

## LINEAR MULTI-CLASS CLASSIFIERS

to predict {hotdog, or not\_hotdog}, a scalar logit  $z$  suffices



for  $K$  classes, use  $K$  logit scores.  
e.g.  $K = 3$ : {hot-dog, pizza, veggie}



for  $K > 2$  classes, a single scalar  $z$  no longer suffices — use  $K$  logit scores to keep track

**softmax:**  $\mathbb{R}^K \rightarrow \mathbb{R}^K$

$$\text{softmax}(z) := \begin{bmatrix} \frac{\exp(z_1)}{\sum_{k=1}^K \exp(z_k)} \\ \vdots \\ \frac{\exp(z_K)}{\sum_{k=1}^K \exp(z_k)} \end{bmatrix}$$

outputs all in  $[0, 1]$ , sum to 1  
non-neg

e.g.

$$\text{softmax} \left( \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} \right) = \begin{bmatrix} \frac{e^1}{e^1 + e^2 + e^3} \\ \frac{e^2}{e^1 + e^2 + e^3} \\ \frac{e^3}{e^1 + e^2 + e^3} \end{bmatrix} = \begin{bmatrix} 0.0900 \\ 0.2447 \\ 0.6653 \end{bmatrix}$$

max among the  $K$  logits

largest # preserved dominance & softer

"soft" max'd in the output

**sigmoid**  $\mathbb{R} \rightarrow \mathbb{R}$

$$\sigma(z) := \frac{1}{1 + \exp(-z)} = \frac{\exp(z)}{\exp(z) + \exp(0)}$$

implicit logit for the negative class

predict positive if  $\sigma(z) > 0.5 = \sigma(0)$

generalization of  $\sigma$   
**softmax:**  $\mathbb{R}^K \rightarrow \mathbb{R}^K$

$$\text{softmax}(z) := \begin{bmatrix} \frac{\exp(z_1)}{\sum_{k=1}^K \exp(z_k)} \\ \vdots \\ \frac{\exp(z_K)}{\sum_{k=1}^K \exp(z_k)} \end{bmatrix}$$

predict the category with the highest softmax score

unifying rule: predict the class with the largest logit

|              | linear logistic binary classifier  | one-out-of- $K$ classifier   |
|--------------|--|--|
| features     | $x \in \mathbb{R}^d$   |  |
| parameters   | $\theta \in \mathbb{R}^d, \theta_0 \in \mathbb{R}$   | $\theta \in \mathbb{R}^{d \times K}, \theta_0 \in \mathbb{R}^K$  |
| linear combo | $\theta^T x + \theta_0 = z \in \mathbb{R}$   | $\theta^T x + \theta_0 = z \in \mathbb{R}^K$   |
| predict      | $\sigma(z) = \frac{\exp(z)}{\exp(0) + \exp(z)}$<br>predict positive if $\sigma(z) > \sigma(0)$ | $\text{softmax}(z) = \begin{bmatrix} \frac{\exp(z_1)}{\sum_{k=1}^K \exp(z_k)} \\ \vdots \\ \frac{\exp(z_K)}{\sum_{k=1}^K \exp(z_k)} \end{bmatrix}$<br>predict the class with the highest softmax score |

## ONE-HOT ENCODING:

- generalizes from {0,1} binary labels
- encode  $K$  classes as  $\mathbb{R}^K$  vector w/ single 1 (hot) & 0s elsewhere

Negative log-likelihood  $K$ -classes loss (aka, cross-entropy)

$g$  : softmax output

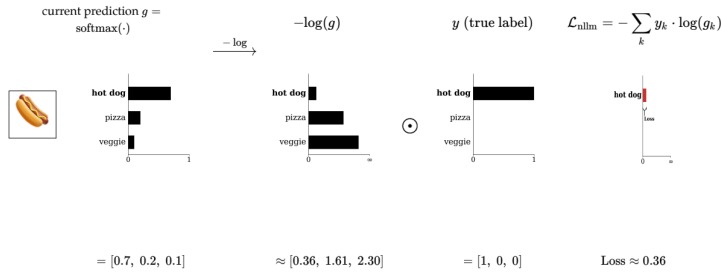
$g_k$  : probability or confidence of belonging in class  $k$

$$\mathcal{L}_{\text{nllm}}(g, y) = - \sum_{k=1}^K y_k \cdot \log(g_k)$$

$y$  : one-hot encoding label

$y_k$  :  $k$ th entry in  $y$ , either 0 or 1

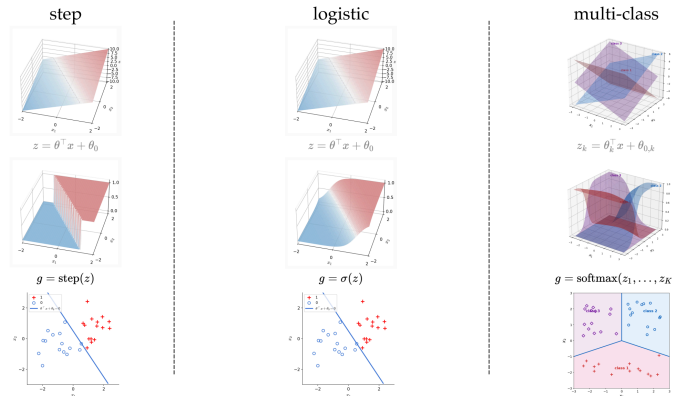
- Generalizes negative log likelihood loss  $\mathcal{L}_{\text{nll}}(g, y) = -[y \log g + (1 - y) \log(1 - g)]$
- Despite the  $K$ -term sum, only the term corresponding to its true class label contributes, since all other  $y_k = 0$



To reduce the loss,  $g_{\text{hot dog}}$  needs to go up — this signal flows smoothly back to  $\theta$  through  $-\log$  and softmax, so we can optimize via gradient descent.

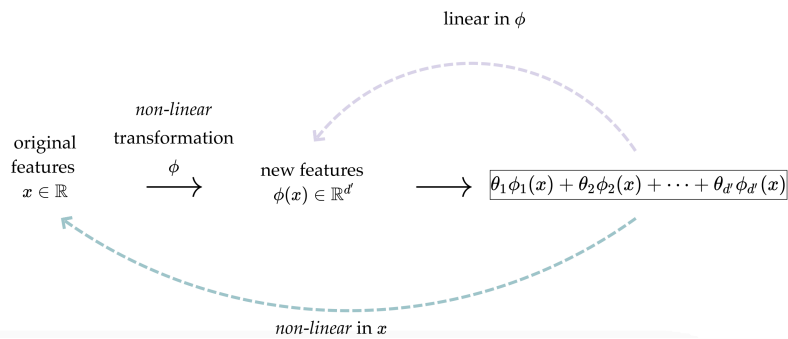
### 3/2 FEATURES, NEURAL NETWORKS

Recall:



decision boundary is linear in feature  $x$

### SYSTEMATIC FEATURE TRANSFORMATIONS



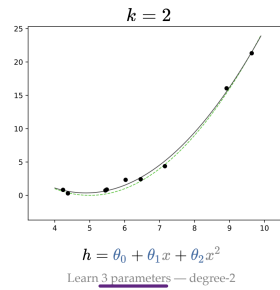
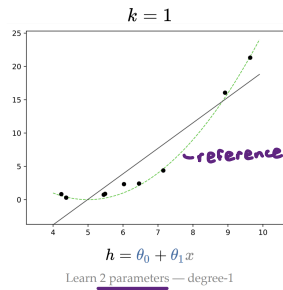
# systematic polynomial features construction

|          | $d = 1$ , features: $x_1$    | $d = 2$ , features: $x_1, x_2$  |
|----------|------------------------------|---|
| $k = 0$  | 1                            | 1   |
| $k = 1$  | 1, $x_1$                     | 1, $x_1, x_2$   |
| $k = 2$  | 1, $x_1$ , $x_1^2$           | 1, $x_1, x_2$ , $x_1^2, x_1x_2, x_2^2$                                      |
| $k = 3$  | 1, $x_1$ , $x_1^2$ , $x_1^3$ | 1, $x_1, x_2$ , $x_1^2, x_1x_2, x_2^2$ , $x_1^3, x_1^2x_2, x_1x_2^2, x_2^3$ |
| $\vdots$ | $\vdots$                     | $\vdots$  |

build up higher & higher terms (more poly relationships)

- Elements in the basis are the monomials of original features raised up to power  $k$
- With a given  $d$  and a fixed  $k$ , the basis is fixed.

- $n = 9$  data points, each with feature  $x \in \mathbb{R}$  and label  $y \in \mathbb{R}$
- data generated from green dashed curve

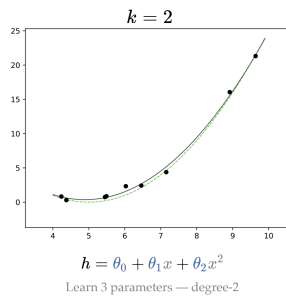
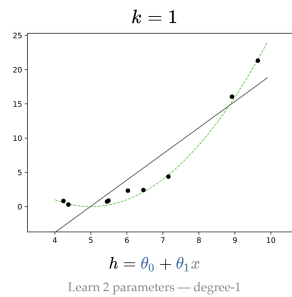


**k**: hyperparameter that determines capacity (expressiveness) of hypothesis class.

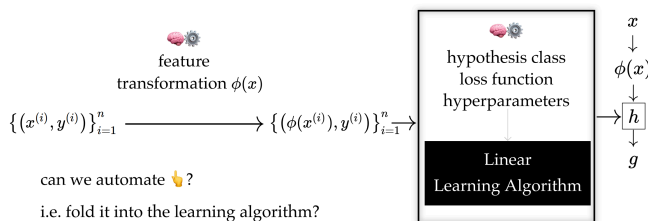
- models w/ many rich features & free params tend to have high capacity but greater risk of overfitting.

- how to choose k? validation/cross validation

- $n = 9$  data points, each with feature  $x \in \mathbb{R}$  and label  $y \in \mathbb{R}$
- data generated from green dashed curve



today, so far:

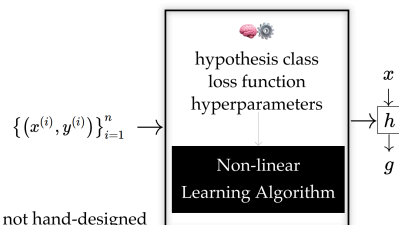


can we automate  $\downarrow$ ?

i.e. fold it into the learning algorithm?

## NEURAL NETWORKS:

neural networks:



the nonlinearity  $\phi$  is learned not hand-designed

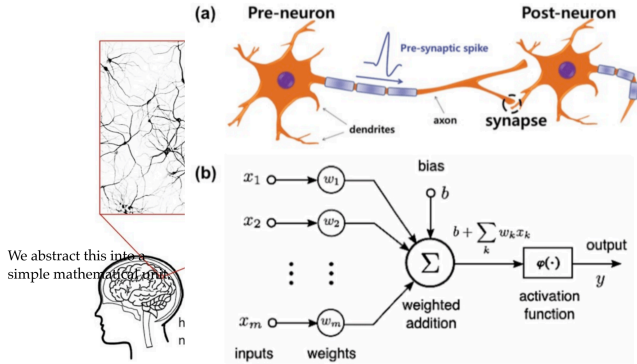
**NN FUNDAMENTALS:**

- nonlinear feature transformation
- composing simple nonlinearities amplifies this effect
- backpropagation ← efficient learning

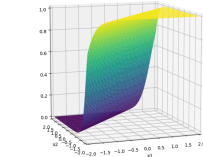
Expressiveness

layered structure

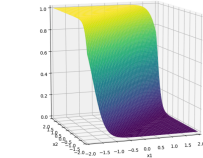
**Neurons and the brain**



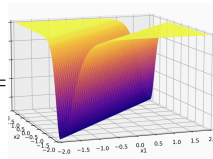
$\sigma_1 = \sigma(5x_1 - 5x_2 + 1)$



$\sigma_2 = \sigma(-5x_1 + 5x_2 + 1)$

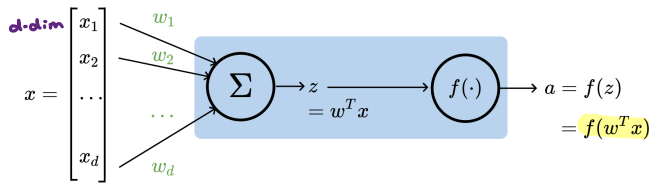


some appropriately weighted sum

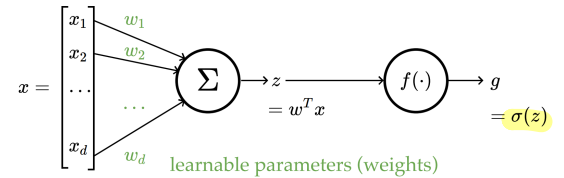


• "Composing" simple transformations

**A neuron:** ☆



- $x$ :  $d$ -dimensional input
- $w$ : weights (i.e. parameters) *old @*  $w$ : what the algorithm learns
- $z$ : pre-activation output *scalar*  $z$ : scalar
- $f$ : activation function *scalar in scalar out*  $f$ : what we engineers choose
- $a$ : post-activation output  $a$ : scalar

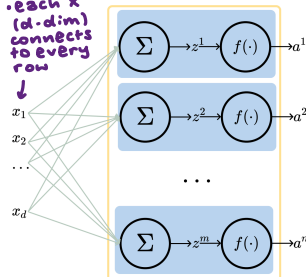


Choose activation  $f(z) = \sigma(z)$

**LAYER:**

A layer:

• each  $x$  ( $d$ -dim) connects to every row

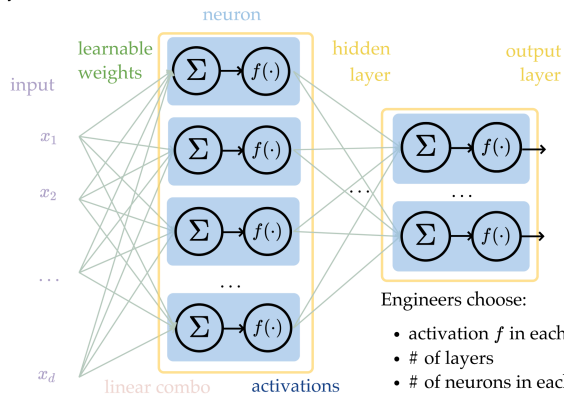


learnable weights

- (# of neurons) = layer's output dimension.
- activations applied element-wise, e.g.  $z^1$  won't affect  $a^2$ . (softmax in the output layer is a common exception)
- all neurons in a layer typically share the same  $f$  (mixed  $f$  is possible but complicates the math).
- layers are typically fully connected, each input  $x_i$  influences every  $a^j$ .

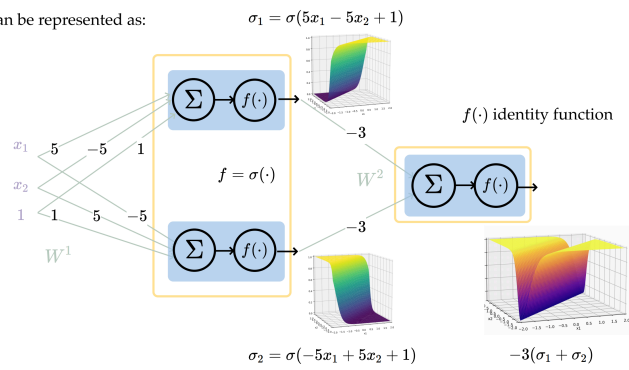
e.g. linear logistic classifier represented as a computation graph

A (fully-connected, feed-forward) neural network:

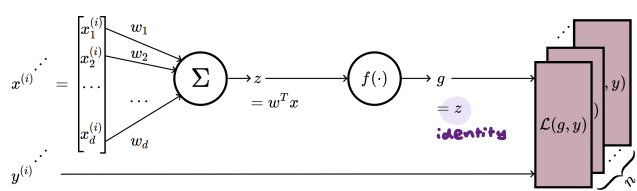


- Engineers choose:
- activation  $f$  in each layer
  - # of layers
  - # of neurons in each layer

can be represented as:

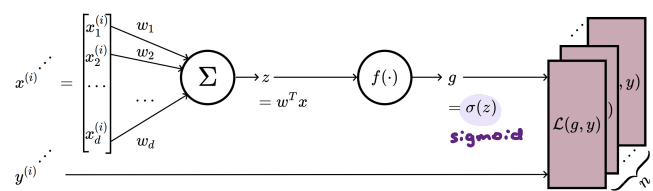


e.g. forward-pass of a linear regressor



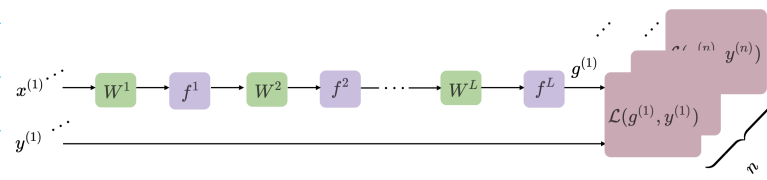
- Activation  $f$  is chosen as the identity function
- Evaluate the loss  $\mathcal{L} = (g^{(i)} - y^{(i)})^2$
- Repeat for each data point, average the sum of  $n$  individual losses

e.g. forward-pass of a linear logistic classifier



- Activation  $f$  is chosen as the sigmoid function
- Evaluate the loss  $\mathcal{L} = -[y^{(i)} \log g^{(i)} + (1 - y^{(i)}) \log (1 - g^{(i)})]$
- Repeat for each data point, average the sum of  $n$  individual losses

Forward pass: evaluate, given the current parameters



- the model outputs  $g^{(i)} = f^L(\dots f^2(f^1(x^{(i)}; \mathbf{W}^1); \mathbf{W}^2); \dots \mathbf{W}^L)$
- the loss incurred on the current data  $\mathcal{L}(g^{(i)}, y^{(i)})$
- the training error  $J = \frac{1}{n} \sum_{i=1}^n \mathcal{L}(g^{(i)}, y^{(i)})$

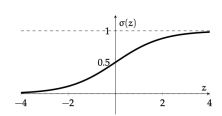
linear combination  
(nonlinear) activation  
loss function

DESIGN CHOICES

Hidden layer activation function  $f$  choices

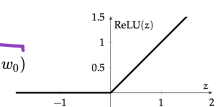
$\sigma$  used to be the most popular

- firing rate of a neuron
- elegant gradient  $\sigma'(z) = \sigma(z) \cdot (1 - \sigma(z))$



nowadays, the default choice:

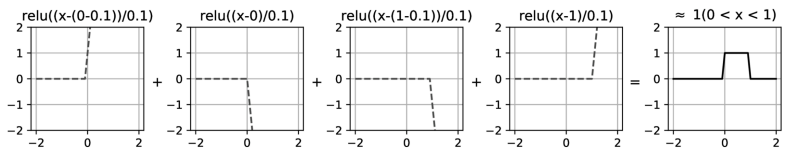
$$\text{ReLU}(z) = \begin{cases} 0 & \text{if } z < 0 \\ z & \text{otherwise} \end{cases} = \max(0, z) = \max(0, w^T x + w_0)$$



very simple function form (so is the gradient).

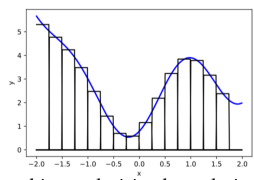
$$\frac{\partial \text{ReLU}(z)}{\partial z} = \begin{cases} 0, & \text{if } z < 0 \\ 1, & \text{otherwise} \end{cases}$$

compositions of ReLU(s) can be quite expressive



compositions of ReLU can be quite expressive (even just 4 ReLU neurons already curves a non-trivial surface)

asymptotically, can approximate any continuous function arbitrarily well (for regression)



therefore can also approximate arbitrary decision boundaries (for classification)

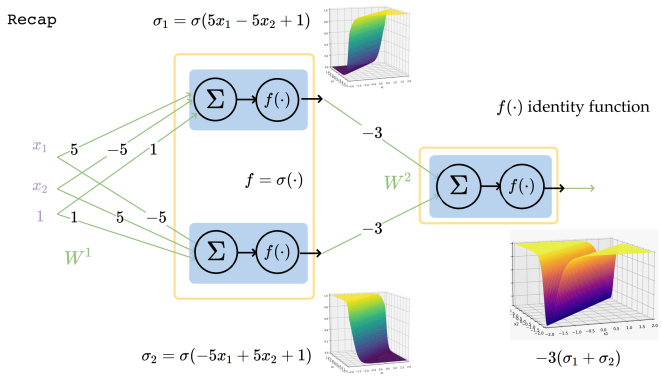
WIDTH: # neurons in layers

DEPTH: # of layers

→ increasing width/depth (w/ non-linear activation) makes model more expressive

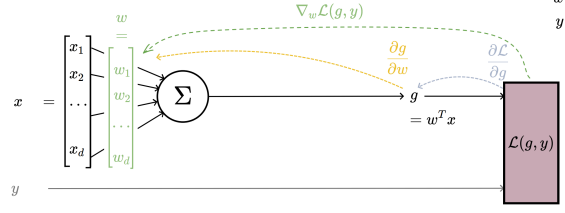
but also incr. risk of overfitting

3/9



**BACKWARD PASS (to learn params / weights)**

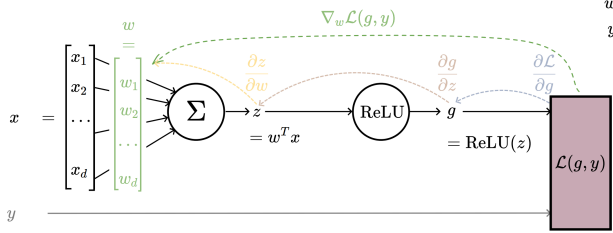
for simplicity, say training data set is just  $(x, y)$  and squared loss



chain rule  

$$\nabla_w \mathcal{L}(g, y) = \frac{\partial \mathcal{L}(g, y)}{\partial w} = \frac{\partial g}{\partial w} \cdot \frac{\partial \mathcal{L}}{\partial g} = \frac{\partial g}{\partial w} \cdot \frac{\partial (g-y)^2}{\partial g} = \frac{\partial (w^T x)}{\partial w} \cdot 2(g-y) = x \cdot 2(g-y)$$

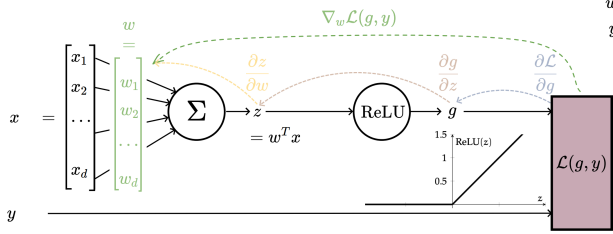
Slightly more interesting:



chain rule  

$$\nabla_w \mathcal{L}(g, y) = \frac{\partial \mathcal{L}(g, y)}{\partial w} = \frac{\partial z}{\partial w} \cdot \frac{\partial g}{\partial z} \cdot \frac{\partial \mathcal{L}}{\partial g} = \frac{\partial z}{\partial w} \cdot \frac{\partial g}{\partial z} \cdot \frac{\partial (g-y)^2}{\partial g} = x \cdot \frac{\partial(\text{ReLU}(z))}{\partial z} \cdot 2(g-y)$$

Slightly more interesting:

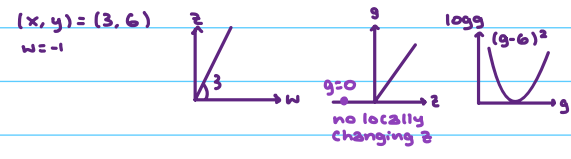
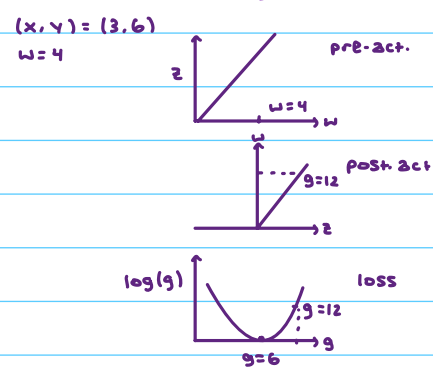
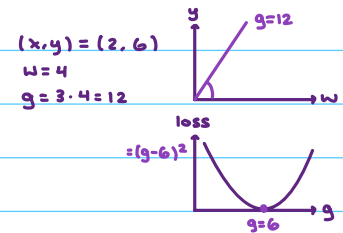


chain rule  

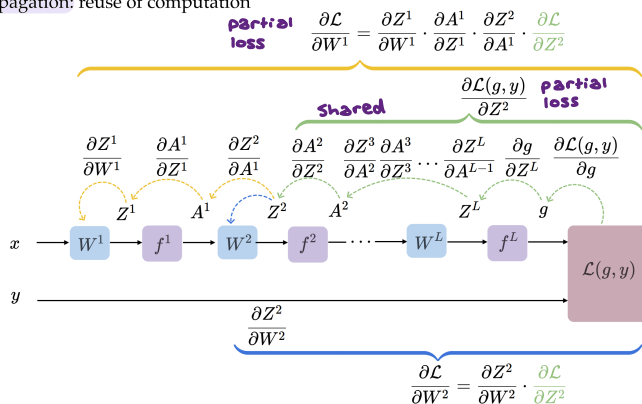
$$\nabla_w \mathcal{L}(g, y) = x \cdot \frac{\partial(\text{ReLU}(z))}{\partial z} \cdot 2(g-y) = \begin{cases} 0, & \text{if } z < 0 \\ 2x(g-y), & \text{otherwise} \end{cases}$$

$$= \begin{cases} 0, & \text{if } z < 0 \\ 1, & \text{otherwise} \end{cases}$$

example on the blackboard



backpropagation: reuse of computation



Training a neural network: the full loop

Initialize all weights  $W^1, W^2, \dots, W^L$  randomly

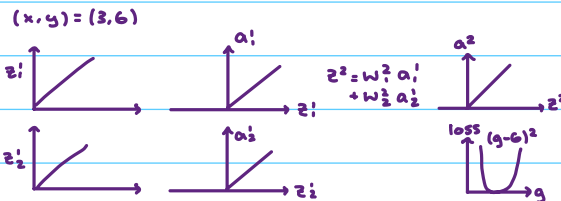
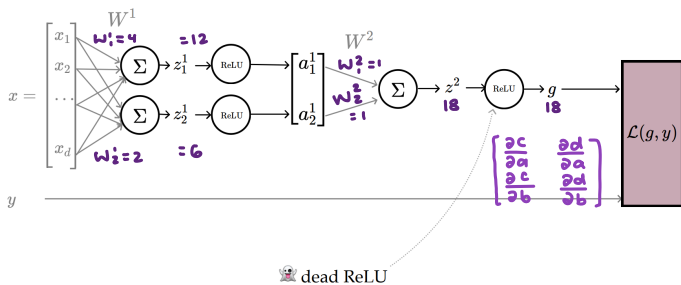
Repeat until stopping criterion:

- Forward pass:** for each data point, compute  $Z^1, A^1, Z^2, A^2, \dots, g^{(i)}$  **evaluate**
- Evaluate loss:** for each data point, compute  $\mathcal{L}(g, y)$
- Backward pass:** pick a data point, compute  $\nabla_{W^\ell} \mathcal{L}(g^{(i)}, y^{(i)})$  for all  $\ell = L, L-1, \dots, 1$  via the chain rule (reuse intermediate results  $\rightarrow$  backpropagation) **to get vals & learn params**
- Update:**  $W^\ell \leftarrow W^\ell - \eta \nabla_{W^\ell} \mathcal{L}$  for all  $\ell$  **\*not guaranteed to be optimal**

#could either evaluate all partial derivatives to numerical vectors or for in chunks

ISSUES & REMEDIES

if  $z^2 < 0$ , no weights get updated

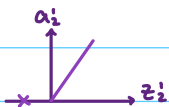


$A \in \mathbb{R}^{a \times a}$ ,  $B \in \mathbb{R}^{b \times b}$

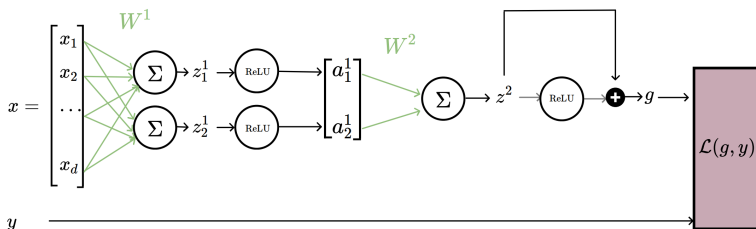
$\frac{\partial A}{\partial B} = [ ] \in \mathbb{R}^{b \times b \times a \times a}$   $\leftarrow (a \times a) \cdot (b \times b)$  partial derivatives to calculate

Ex:  $y = -1, z_2 = -3$

input into ReLU is 0



Residual (skip) connection:



Now,  $g = z^2 + \text{ReLU}(z^2)$   
even if  $z^2 < 0$ , with skip connection, weights in earlier layers can still get updated

# Vanishing and Exploding Gradients

Backpropagation computes gradients via the chain rule — a product of many factors:

$$\text{e.g. } \frac{\partial \mathcal{L}}{\partial W^1} = \frac{\partial Z^1}{\partial W^1} \cdot \underbrace{\frac{\partial A^1}{\partial Z^1} \cdot \frac{\partial Z^2}{\partial A^1} \cdots \frac{\partial A^{L-1}}{\partial Z^{L-1}} \cdot \frac{\partial Z^L}{\partial A^{L-1}}}_{L-1 \text{ layers of multiplicative factors}} \cdot \frac{\partial g}{\partial Z^L} \cdot \frac{\partial \mathcal{L}}{\partial g}$$

If each factor is small, their product shrinks very quickly with depth, little learning.

If factors are large (magnitude > 1), gradients explode, also problematic.

Many practical remedies: residual connection, gradient clipping

Gradient clipping: If  $\|\nabla\| > \tau$ , rescale:  $\nabla \leftarrow \tau \frac{\nabla}{\|\nabla\|}$ .

Preserves gradient direction, caps its magnitude

## Summary

- Multi-layer perceptrons automatically learn good features and transformations from data.
- Training loop: forward pass  $\rightarrow$  loss  $\rightarrow$  backward pass  $\rightarrow$  weight update, repeat.
- Backpropagation reuses intermediate computations to efficiently evaluate all gradients via the chain rule.
- Dead ReLU neurons (pre-activation always negative) get zero gradient and stop learning; careful initialization helps.
- Vanishing/exploding gradients arise from multiplying many small (or large) factors across layers.

---

Reference: weight initialization

- If all weights start at 0, every neuron computes the same thing  $\rightarrow$  no learning
- If weights are too large, activations saturate or explode
- If weights are too small, signals shrink to zero across layers

**Goal:** keep the variance of activations roughly constant across layers.

- **Xavier initialization** (for sigmoid / tanh):  $W_{ij}^\ell \sim \mathcal{N}\left(0, \frac{1}{n_{\ell-1}}\right)$

where  $n_{\ell-1}$  = number of inputs (fan-in) to the layer

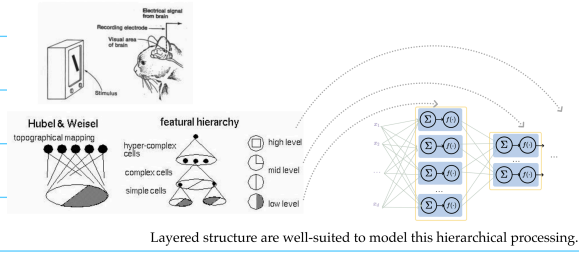
- **He initialization** (for ReLU):  $W_{ij}^\ell \sim \mathcal{N}\left(0, \frac{2}{n_{\ell-1}}\right)$

factor of 2 compensates for ReLU zeroing out half the inputs

3116 CONVOLUTIONAL NEURAL NETWORKS

- 784 inputs → must calculate 13k params
- have a specialized, structured network to help fight overfitting & partly fully-connected nets don't scale well for vision tasks

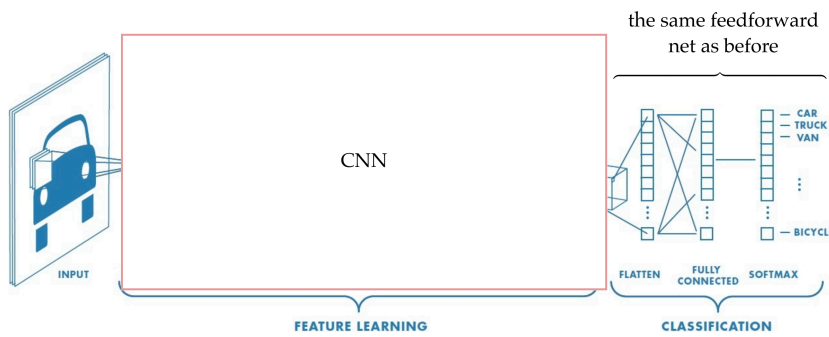
• Visual hierarchy



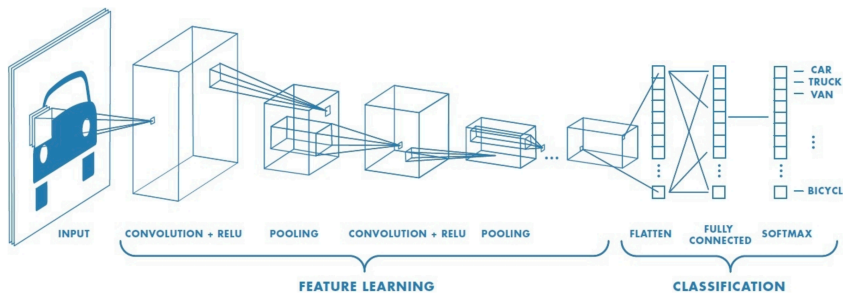
- visual hierarchy
- spatial locality (zoom in & compare things around a pt)
- translational invariance (should recognize all types despite differences)

\* CNN Exploits via layered structure, convolution, and pooling to handle images efficiently & effectively

typical CNN architecture for image classification



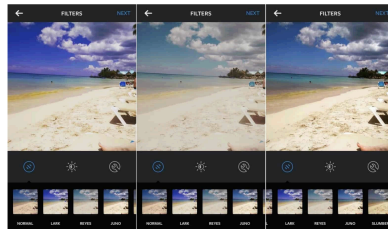
typical CNN structure for image classification



Convolutional layer might sound foreign, but it's very similar to a fully-connected layer

| Layer           | Forward pass, <i>do</i> | Backward pass, <i>learn</i> | Design choices     |
|-----------------|-------------------------|-----------------------------|--------------------|
| fully-connected | dot-product, activation | neuron weights              | neuron count, etc. |
| convolutional   | convolution, activation | filter weights              | conv specs, etc.   |

Convolution result:



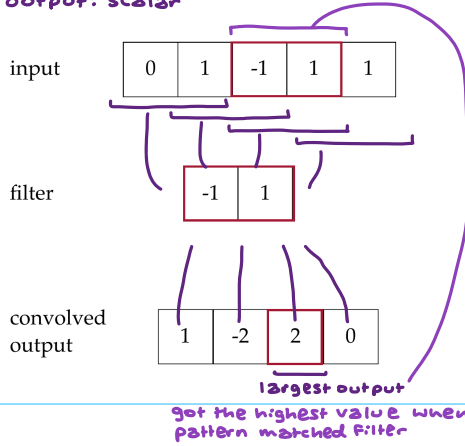
**CONVOLUTION:**

→ template matching (the filter): 1D

convolution interpretation 1:

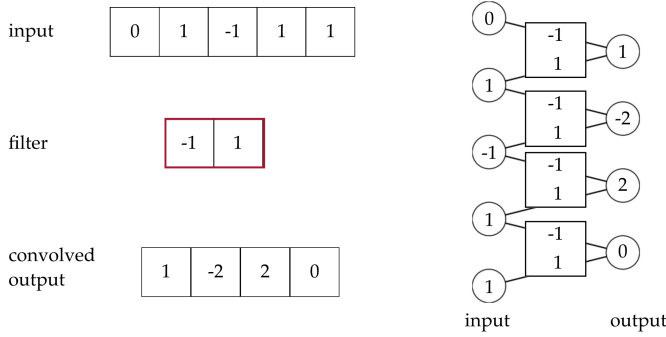
template matching

sliding window, multiply  
output: scalar



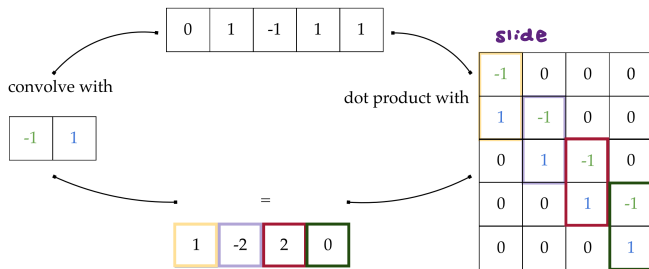
convolution interpretation 2:

"look" locally through the filter  
this local region = receptive field

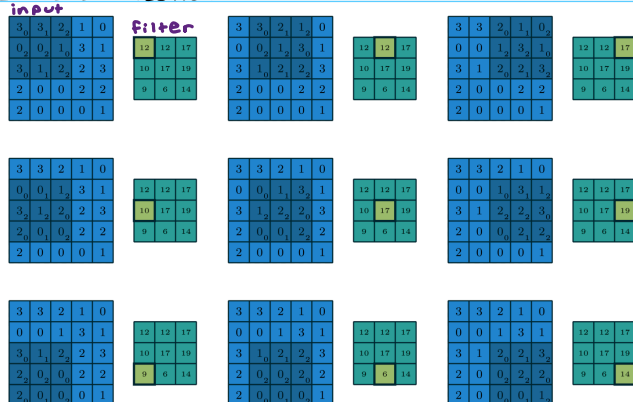


convolution interpretation 3:

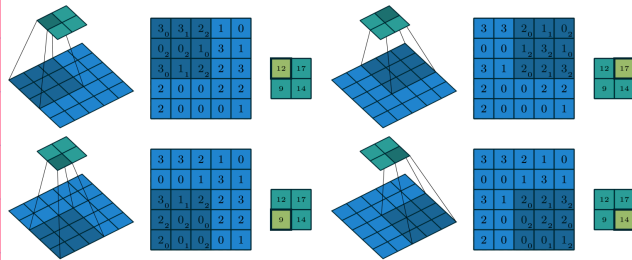
sparse-connected layer  
with parameter sharing



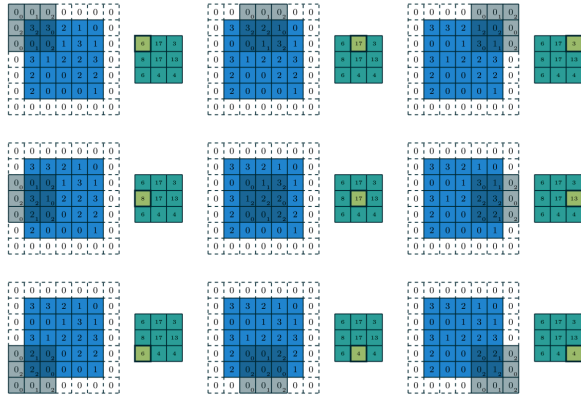
**EX: 2D CONVOLUTION**



**STRIDE OF 2:**

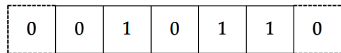


**STRIDE OF 2, PADDING SIZE 1**

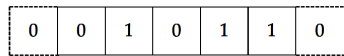


quick summary: hyperparameters for 1d convolution

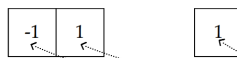
- Zero-padding



- Stride (e.g. stride of 2)



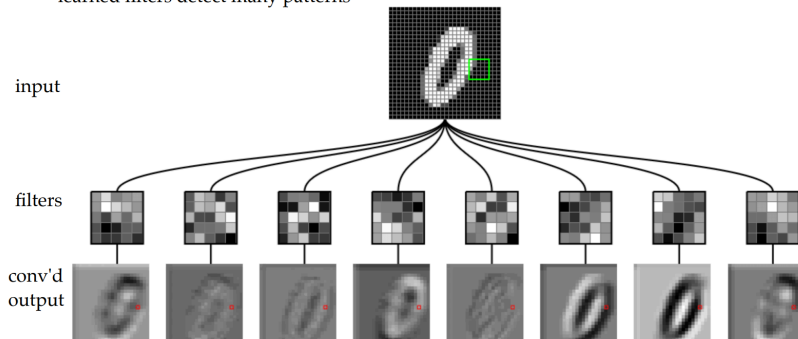
- Filter size (e.g. we saw these two in 1-d)



these weights are what CNN learn eventually

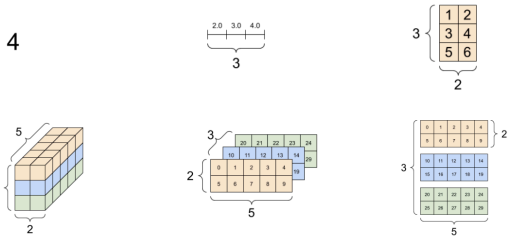
**Each filter helps w/ pattern recognition:**

learned filters detect many patterns



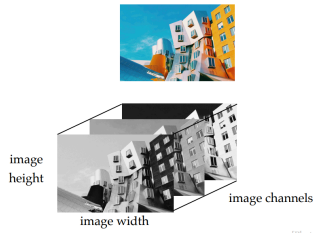
### 3D TENSORS

3D tensors:  
2D matrices  
stacked tog.

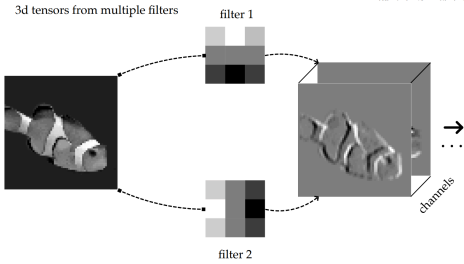


• collection of 2D images (ex: R, G, B) → each channel is a complete but independent view of the same scene.

3d tensors from color channels



3d tensors from multiple filters



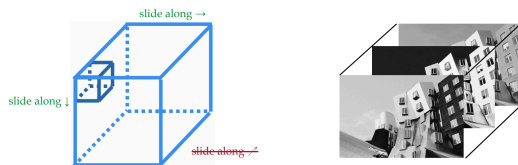
can stack complete but ind. channels

### CHANNELS come from:

1) color input (width x height channels)

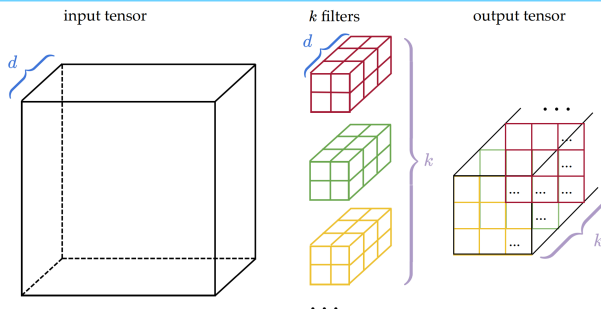
2) multiple filters → multiple channels

Why we don't typically do 3d convolution



Convolution shares weights across shifted positions  
shifting makes sense spatially (a cat can be anywhere)  
but not across channels (red ≠ shifted green)

3D conv is used when the third axis is spatial/temporal (MRI, video)

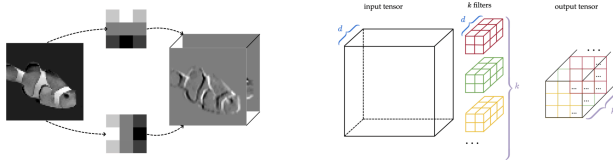


Every convolutional layer works with 3d tensors:

1. color input



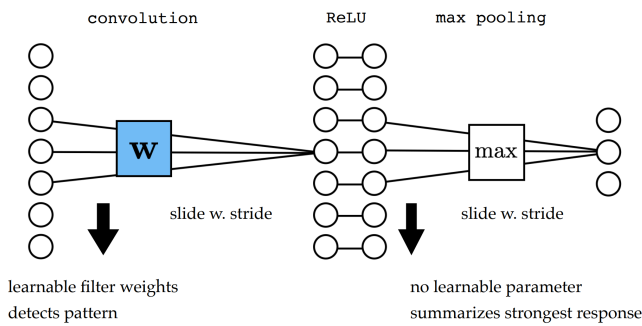
2. the use of multiple filters in doing 2d convolution



**MAX POOLING: & RELU** → 'side where the pattern was; i just care if it existed'

**POOLING:** similar to convolution, but instead of learning a filter, use max pooling filter to

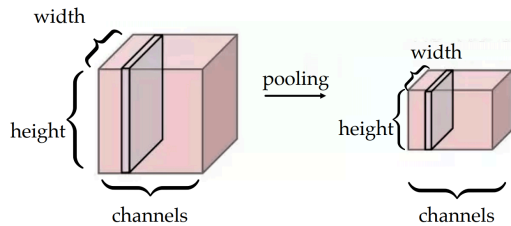
1d max pooling



\*pooling across spatial locations achieves invariance w.r.t. small translations

- large response regardless of exact position of edge

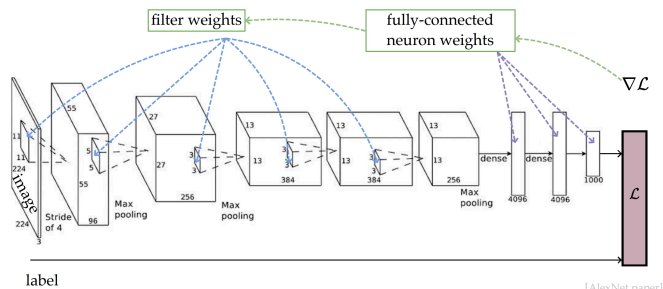
Pooling across *spatial* locations achieves invariance w.r.t. small translations:



applied independently across all channels

so the *channel* dimension remains *unchanged*

**imagenet classification error**



[AlexNet paper]  
[all max pooling are 3x3 filter, stride 2; pooled outputs not explicitly shown on diagram: 27x27x96, 13x13x256, 6x6x256]

Here are the building blocks we will use to build a CNN:

- Convolutional layers
- Max pooling layers
- Fully connected layers

Assume that the task of our CNN is to detect animals like cats, dogs, lizards etc. in an image.

### 2.1)

Which type of layer would be best to use to detect low-level features (like lines, curves, etc.) of an animal in an image?

Choose one:  'pattern matcher'

Ask for Help 100.00%

You have infinitely many submissions remaining.

**Solution:** Convolutional

**Explanation:**  
Filters in the first convolutional layers are responsible for detecting low-level features (e.g., edges, color, contrast). Later convolutional layers are responsible for detecting mid-level features (e.g., ears, eyes).

designed to look for specific local structure.  
(1 surrounded by 0s)

### 2.2)

Suppose we are given a set of inputs where each input represents whether a useful feature of animals (tail, beak, whiskers, etc.) occurs anywhere in a image. Which type of layer would be best to use to turn these inputs into a classification of whether the image is a cat, dog, lizard, etc.?

Choose one:  OR gate: the 'decision maker'

Ask for Help 100.00%

You have infinitely many submissions remaining.

**Solution:** Fully connected

**Explanation:**  
Fully connected layers allow combining features from the entire image and provide the final network output.

1 if image contains at least one instance  
of the filter (w/  $\max(x, 0)$  applied to it)

### 2.3)

Suppose you have detected occurrences of some useful low-level feature, like a furry tail, throughout your image. You would like to use these local detections to figure out whether or not the furry tail occurred anywhere within larger regions of the image. Which layer would be best to use?

Choose one:  'evidence gatherer'

Ask for Help 100.00%

You have infinitely many submissions remaining.

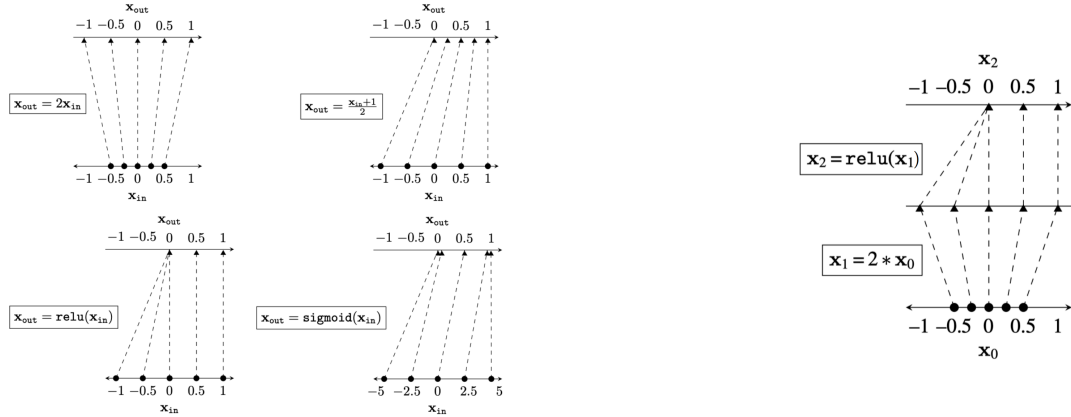
**Solution:** Max pool

**Explanation:**  
Max pooling layers detect the strongest response within a given window. This property allows the network to be less sensitive to feature locations.

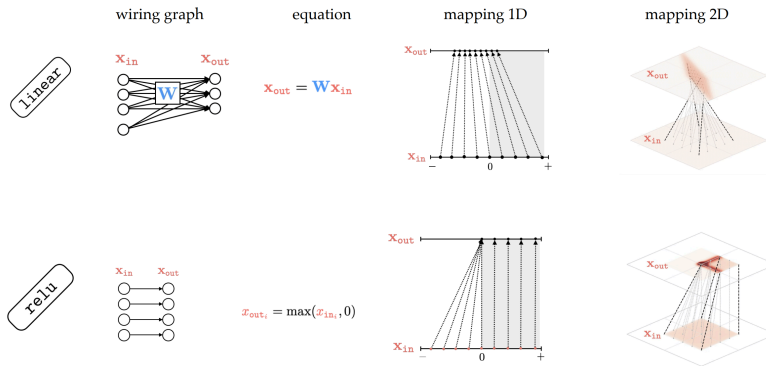
$f(x) = \max(x, 0)$   
ReLU: filter for relevance (keep matches but  
turns noise to 0)  
max pooling: summarizes the region.  
'ide where the pattern was;  
i just care if it existed'

### 3/30 REPRESENTATION LAYERS:

Representation transformations for a variety of neural net operations and stack of neural net operations



#### Parameters



### TRAINING NN CLASSIFIER:

$$g = \text{softmax}(z_2)$$

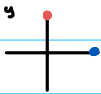
$$z_2 = \text{linear}(a_1) \in \mathbb{R}^2$$

$$a_1 = \text{ReLU}(z_1)$$

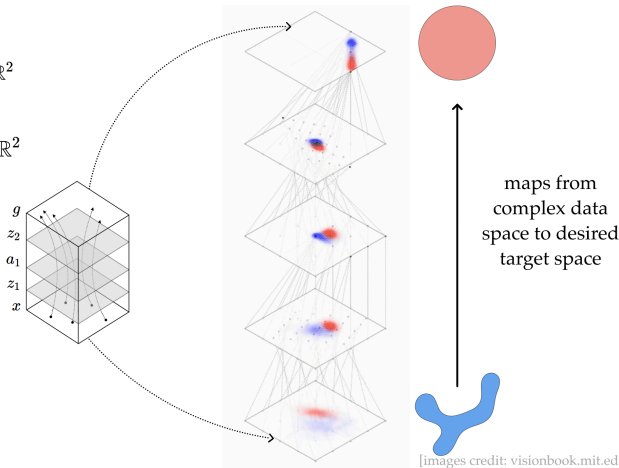
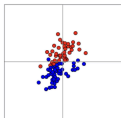
$$z_1 = \text{linear}(x) \in \mathbb{R}^2$$

$$x \in \mathbb{R}^2$$

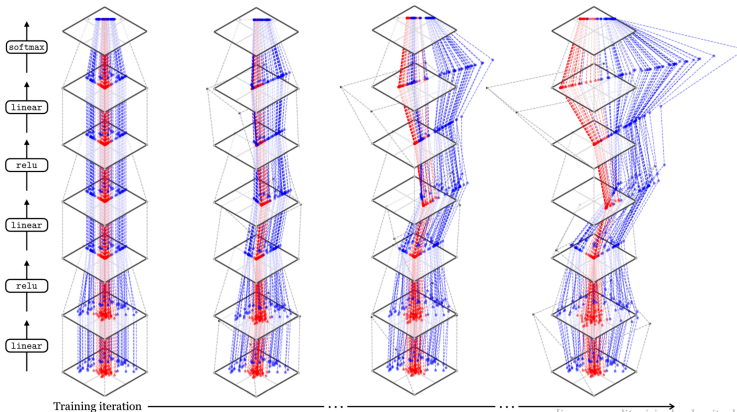
target out:



Training data



[images credit: visionbook.mit.edu]<sub>10</sub>



Training iteration

[images credit: visionbook.mit.edu]<sub>11</sub>

- learnable weights help w/ different scattering
- each layer = different representation (embedding) of data
- from latent embeddings to data: gen. modelling

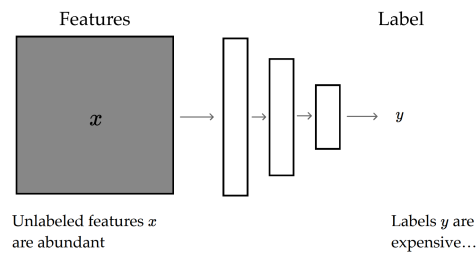
humans also learn representations:

- compact (minimal)
- explanatory (roughly sufficient)
- disentangled (independent factors)
- interpretable (understandable)
- make subsequent problem solving easy

### SELF-SUPERVISED LEARNING

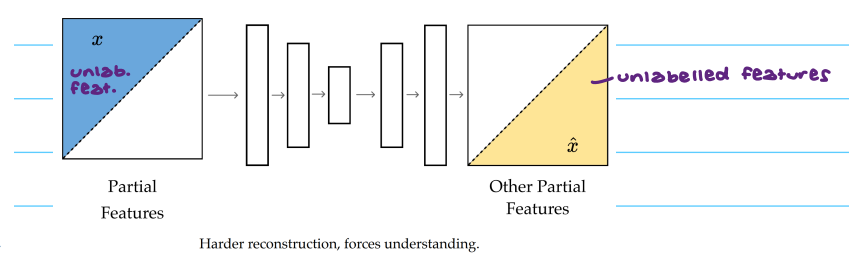
Supervised:

Label prediction (supervised learning)



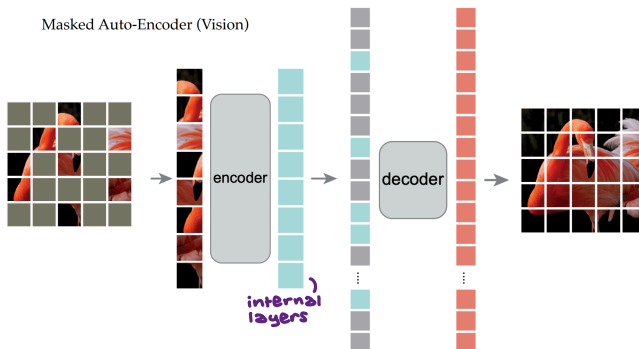
Self-supervised:

Self-supervised Learning (partial feature reconstruction)



lots of unlabelled features on the internet!

Masked Auto-Encoder (Vision)



MASKING: predict color from grayscale (one method)

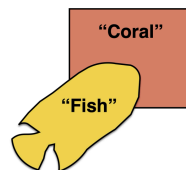
w/o relying on labels

[He, et al. Masked Autoencoders Are Scalable Vision Learners, 2021]

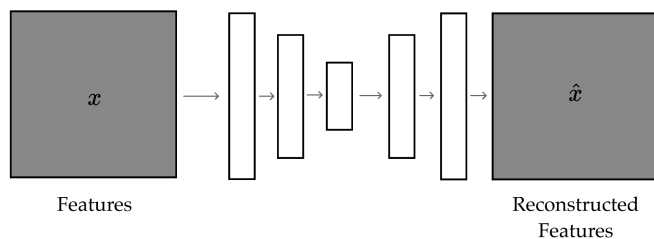
### AUTO-ENCODERS

Good representations are:

- Auto-encoders explicitly aims
- Compact (minimal)
  - Explanatory (roughly sufficient)
  - Disentangled (independent factors)
- these may just emerge as well
- Interpretable (understandable by humans)
  - Make subsequent problem solving easy

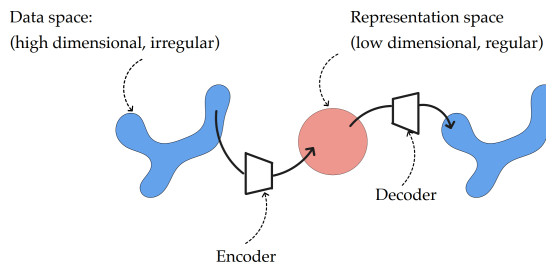
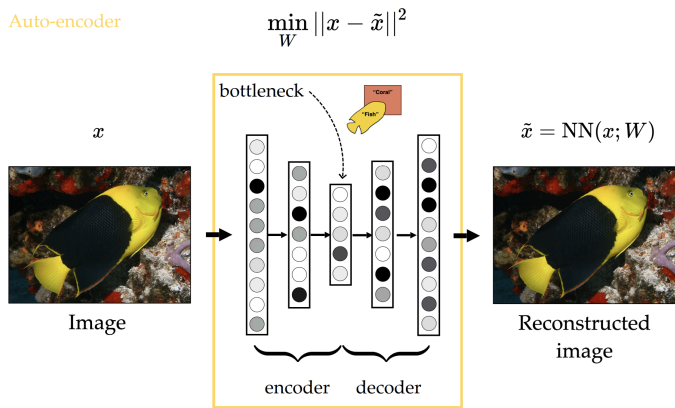


### UNSUPERVISED LEARNING (FEATURE RECONSTRUCTION)

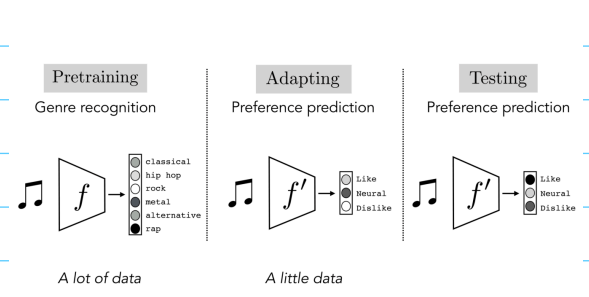
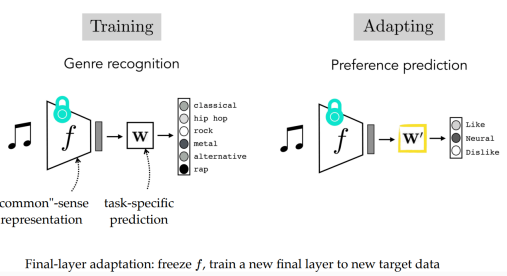
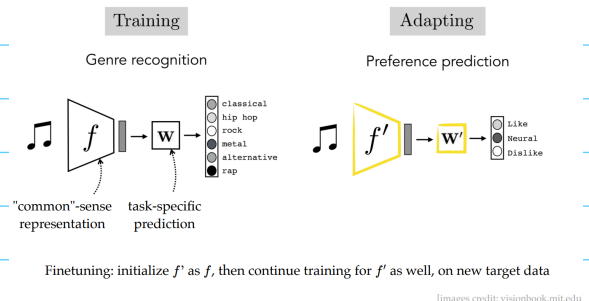
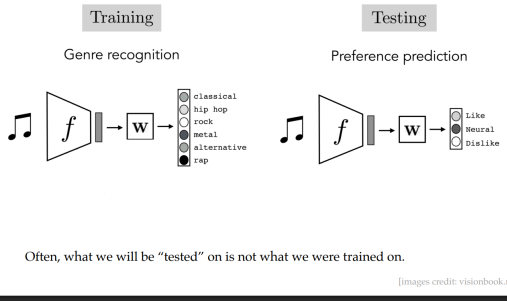


**AUTO-ENCODER:**

Auto-encoder

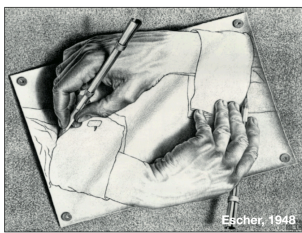


**training & testing:**



**WORD EMBEDDINGS:**

Large Language Models (LLMs) are trained in this self-supervised way

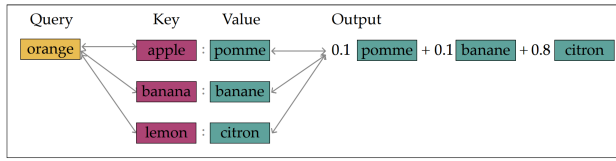


embed words → vectors

(2 words close semantically will be closer together)

- Scrape the internet for plain texts.
- Cook up "labels" (prediction targets) from these texts.
- Convert "unsupervised" problem into "supervised" setup.

\*good word embeddings space is equipped w/ sensible dot-prod. similarity



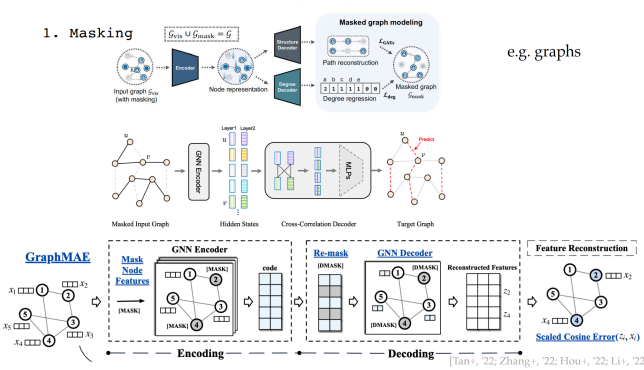
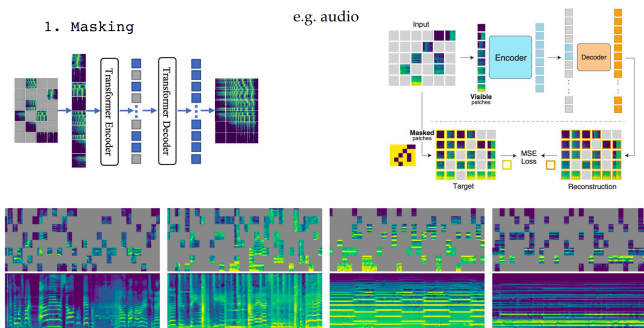
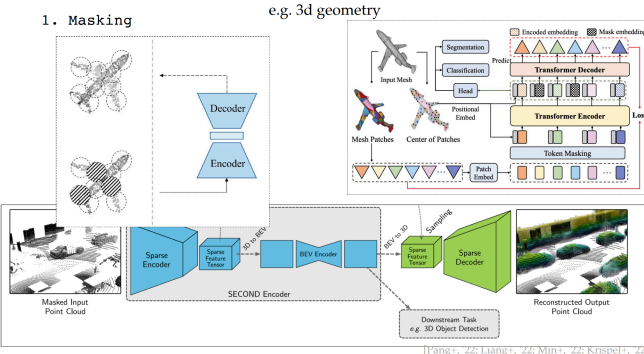
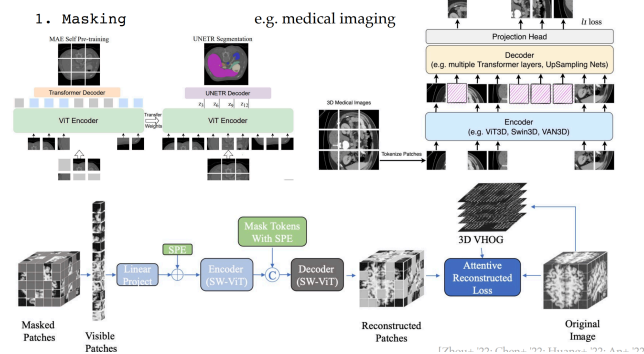
We put (query, key, value) in "good" embeddings in our human brain

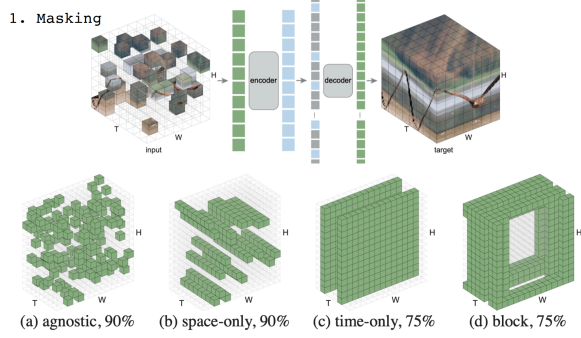
such that "merging" the values  $0.1 \text{ pomme} + 0.1 \text{ banane} + 0.8 \text{ citron}$

via these merging percentages [0.1 0.1 0.8] made sense

very roughly, the attention mechanism in transformers automates this process. We'll see the details next week.

**MASKING:**





\*masking 75% is optimal for images  
 15% optimal for languages  
 95% optimal for videos

[Feichtenhofer, et al., "Masked Autoencoders As Spatiotemporal Learners", NeurIPS 2022] 82

**CONTRASTIVE LEARNING:**

- glimpse of an object might have association w/ smthg else
- crop parts of image & construct network s.t. patches coming from same image are encoded close together

**MULTI-MODALITY**

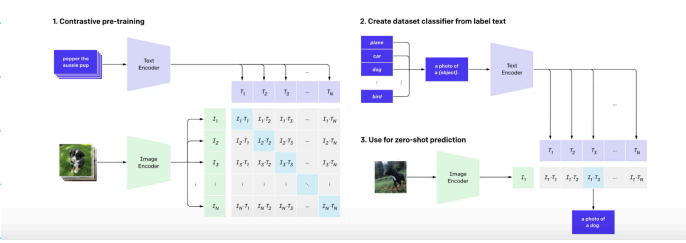
- predict sound cow makes based on image of cow
- ex: input = only video sound of ocean, LLM output = sound of ocean

e.g. image classification (done in the contrastive way)



[Radford et al, Learning Transferable Visual Models From Natural Language Supervision, ICML, 2011] 83

e.g. image classification (done in the contrastive way)

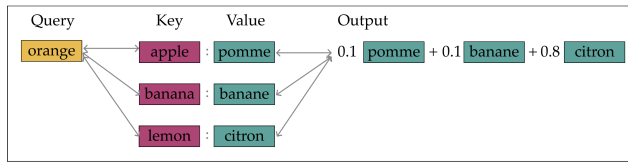


**How Much Information is the Machine Given during Learning?** Y. LeCun

- ▶ **"Pure" Reinforcement Learning (cherry)**
  - ▶ The machine predicts a scalar reward given once in a while.
  - ▶ **A few bits for some samples**
- ▶ **Supervised Learning (icing)**
  - ▶ The machine predicts a category or a few numbers for each input
  - ▶ Predicting human-supplied data
  - ▶ **10→10,000 bits per sample**
- ▶ **Self-Supervised Learning (cake génoise)**
  - ▶ The machine predicts any part of its input for any observed part.
  - ▶ Predicts future frames in videos
  - ▶ **Millions of bits per sample**



4/6 recap: word embedding



a. compare query and key for merging percentages:

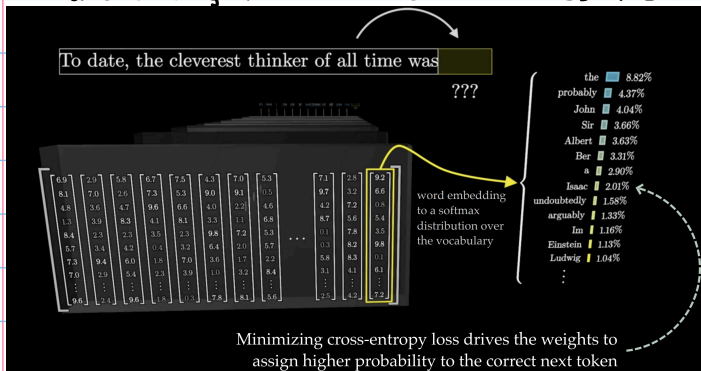
$$\text{softmax} \left( \begin{matrix} \text{orange} & \text{apple} \\ \text{orange} & \text{banana} \\ \text{orange} & \text{lemon} \end{matrix} \right) = [0.1 \ 0.1 \ 0.8]$$

b. then output mixed values 0.1 pomme + 0.1 banane + 0.8 citron

Let's see how this intuition becomes a trainable mechanism.

TRANSFORMERS:

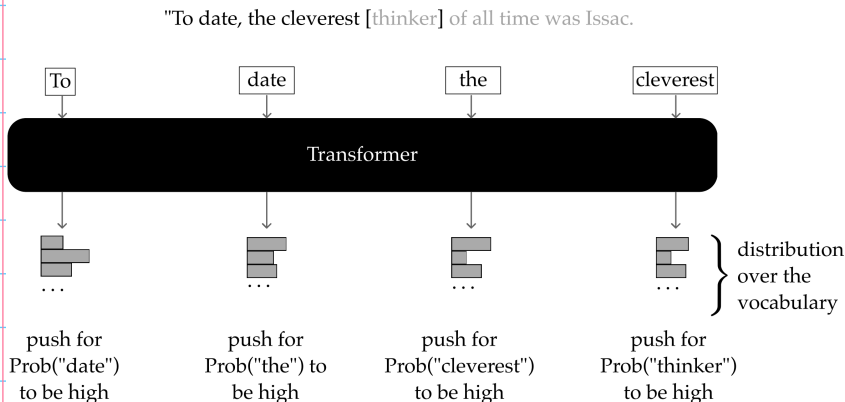
- input embedding: n small chunks, each of which w/ a mapping (ex: GPT-2 n=768, GPT-3 n=1200+)
- passed through attention layer (values internally change) then multi-layer perceptron
- word embedding to softmax dist. over vocab (want correct word to be assigned high prob.)

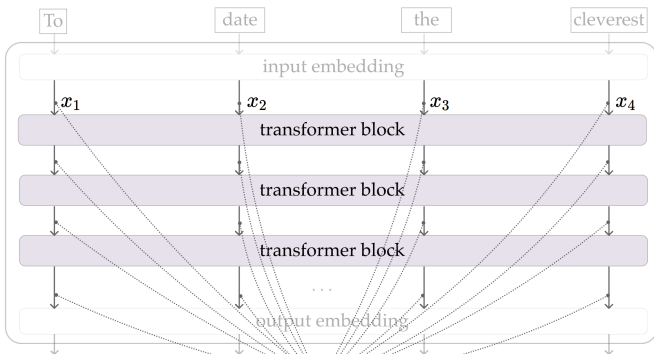
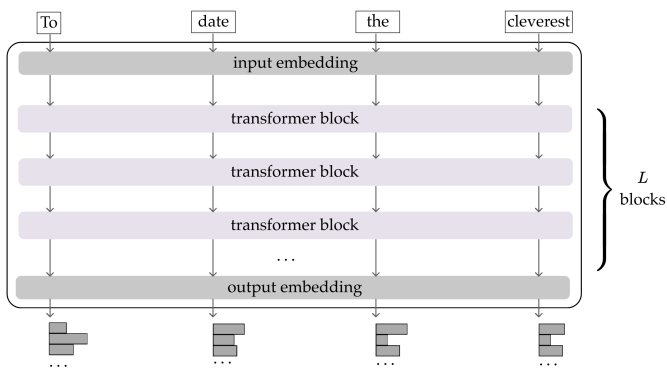


Minimizing cross-entropy loss drives the weights to assign higher probability to the correct next token

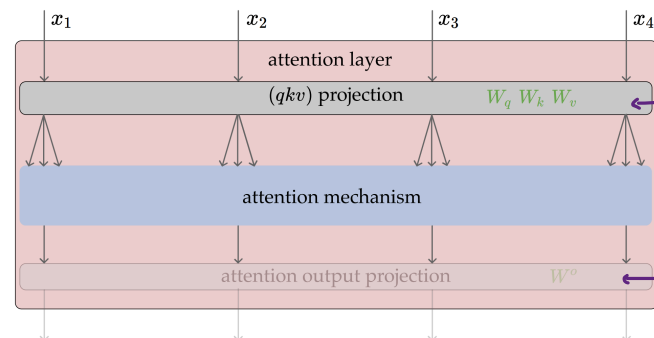
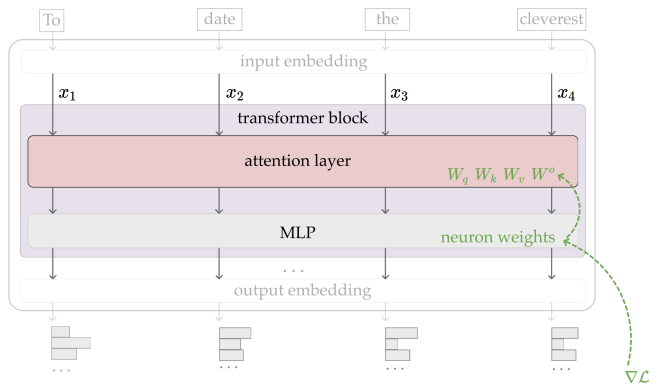
WORDS W/ MULTIPLE MEANINGS:

- all 3 "moles" mapped to same representation E(mole)
- look @ surrounding words & soak in context to get more specific representation





each of the  $n$  tokens transformed, block by block within a shared  $d$ -dimensional word-embedding space.



4 sets of weights we can learn ( $W_q, W_k, W_v, W^o$ )

Most important bits in an attention layer:

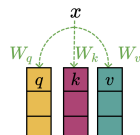
1. (query, key, value) projection
2. attention mechanism

### 1. (query, key, value) projection: linear transformation

1. (query, key, value) projection

With learned projections, we frame  $x$  into:

- a query to be the questions
- a key to be compared
- a value to contribute



Why learning these projections:

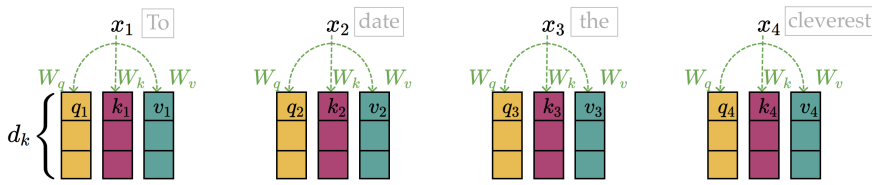
- $W_q$  learns how to ask → what to ask
- $W_k$  learns how to listen
- $W_v$  learns how to speak → giving an answer

$d$ : word embedding

$d_k$ :  $(qkv)$ -embedding ← learned way of phrasing good value/response

attention head:

### 1. (query, key, value) projection

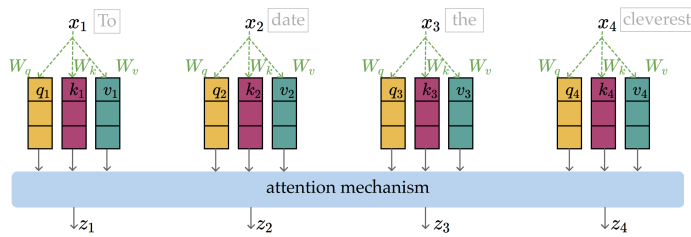


- $W_q, W_k, W_v$ , all in  $\mathbb{R}^{d \times d_k}$
- project word embedding  $x$  from  $d$ -dimensional space to  $d_k$ -dimensional ( $q, k, v$ ) spaces (typically  $d_k < d$ )
- $q_i = W_q^T x_i, k_i = W_k^T x_i, v_i = W_v^T x_i, \forall i$  — *weight sharing* across positions
- *parallel* and *structurally identical* processing

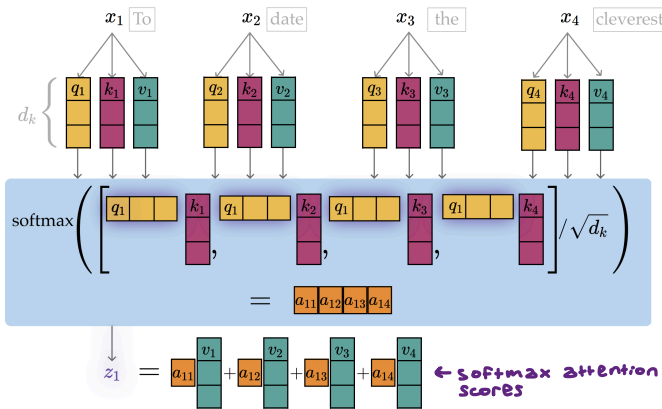
⊛ ONCE YOU UNDERSTAND HOW TO PHRASE GOOD QUESTION, IT'S VERY REUSEABLE

### 2: attention mech.

#### 2. Attention mechanism



- Attention mechanism turns the projected ( $q, k, v$ ) into  $z$
- Each  $z$  is context-aware: a mixture of everyone's values, weighted by relevance

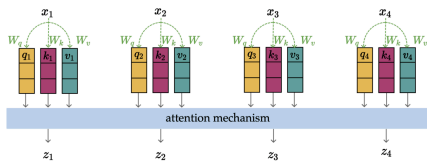


PARALLEL & STRUCTURALLY IDENTICAL PROCESSING: can calculate  $z_4$  w/o  $z_3$

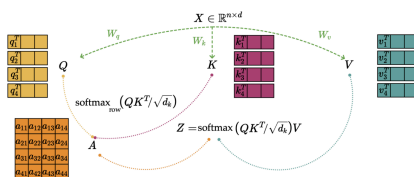
$W_q, W_k, W_v$  shared w/ all elements! but even though question meaning is the same, way question is asked might be slightly different

### ATTENTION HEAD: compact matrix form

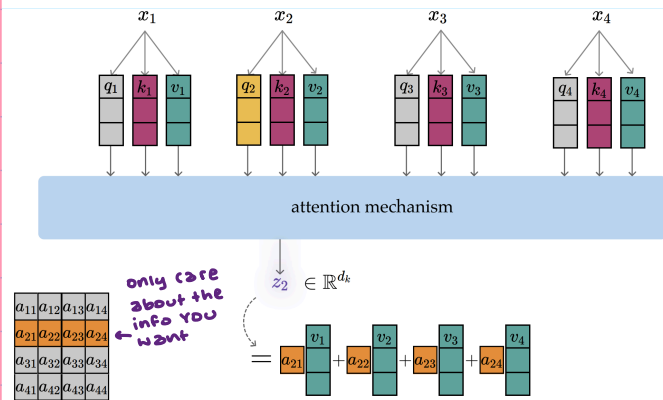
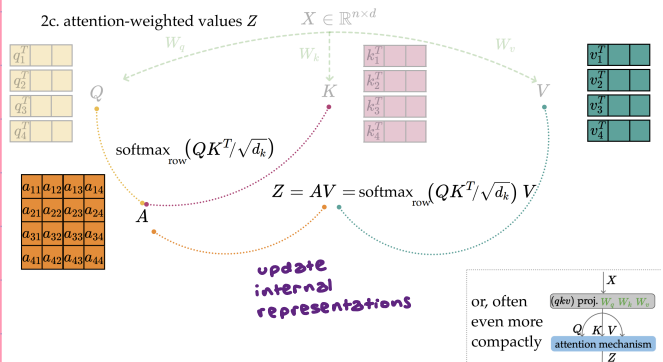
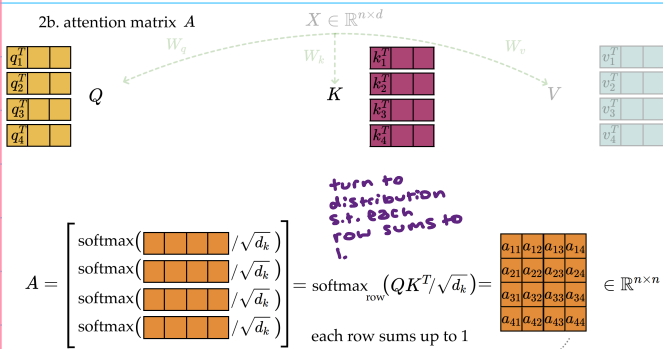
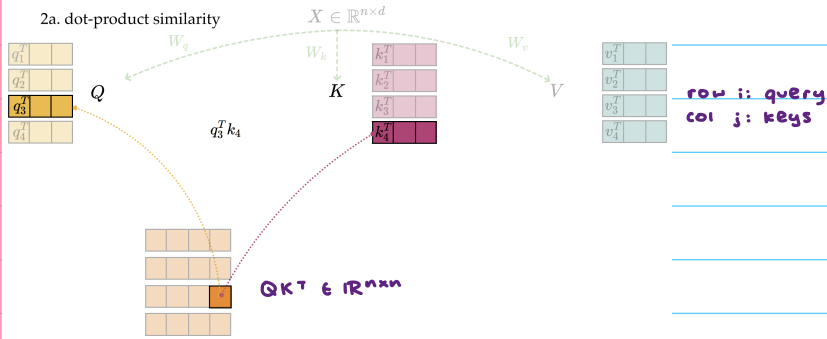
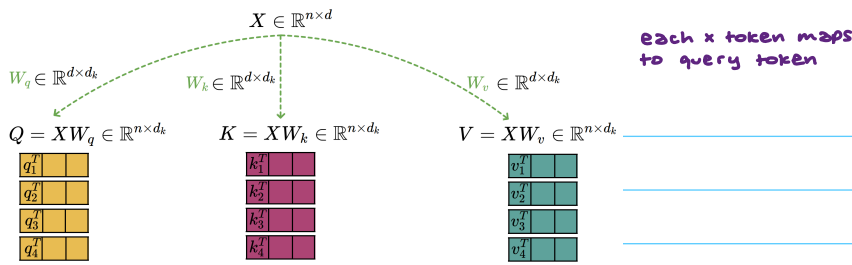
#### Attention head - compact matrix form



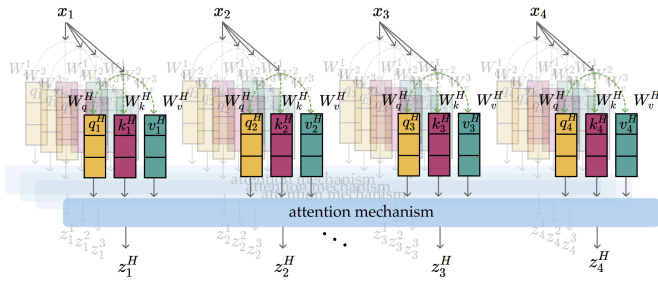
⇕ By stacking each individual vector in the sequence as a row



1. (query, key, value) projection



**MULTI-HEAD ATTENTION:**

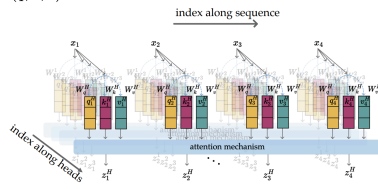


Multi-head Attention

Parallel, and structurally identical processing across all heads and tokens.

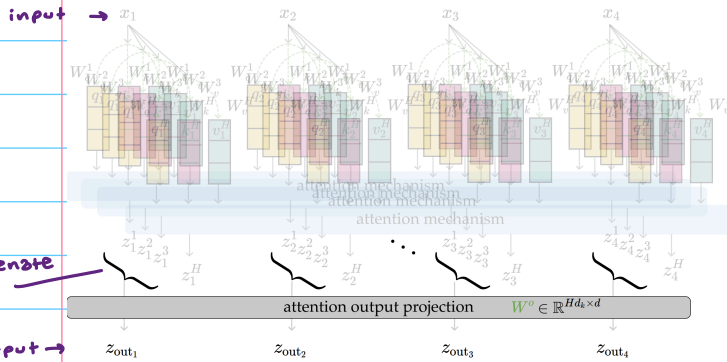
In particular, each head:

- learns its own set of  $W_q, W_k, W_v$
- creates its own projected sequence of  $(q, k, v)$
- computes its own sequence of  $z$
- structurally identical processing
- for each token in the sequence:
  - structurally identical processing



\* big parallel-processing machine

↳ you can learn @ diff. paces but still all learn from each other's questions (parallel- don't have to wait for others to learn)



concatenate as z.

output →



Shape Example:

$n$  num tokens 5  
 $d$  word-embedding dim 6  
 $d_k$  (qkv) embedding dim 3

for a single attention head

|         |               |                |              |
|---------|---------------|----------------|--------------|
| $X$     | input         | $n \times d$   | $5 \times 6$ |
| $W_q^h$ | query proj    | $d \times d_k$ | $6 \times 3$ |
| $W_k^h$ | key proj      | $d \times d_k$ | $6 \times 3$ |
| $W_v^h$ | value proj    | $d \times d_k$ | $6 \times 3$ |
| $Q^h$   | query         | $n \times d_k$ | $5 \times 3$ |
| $K^h$   | key           | $n \times d_k$ | $5 \times 3$ |
| $V^h$   | value         | $n \times d_k$ | $5 \times 3$ |
| $A^h$   | attn matrix   | $n \times n$   | $5 \times 5$ |
| $Z^h$   | attn head out | $n \times d_k$ | $5 \times 3$ |

$H$  num heads 2

concat( $Z^1 \dots Z^H$ )

|           |                |                 |              |
|-----------|----------------|-----------------|--------------|
|           | multi-head out | $n \times Hd_k$ | $5 \times 6$ |
| $W^o$     | output proj    | $Hd_k \times d$ | $6 \times 6$ |
| $Z_{out}$ | attn layer out | $n \times d$    | $5 \times 6$ |

$W$ s are the learned weights

hyperparams →

LECTURE 10: MARKOV DECISION PROCESS ← this + LL create reinforcement learning

POLICY EVAL: similar to fwd PASS; given smthg & evaluating quality of that thing

POLICY OPT: similar to back PASS; searching for smthg of 'good quality'

\* finding good policy similar to finding good set of weights

S: state space, contains all possible states  $S$

A: action space, contains all possible actions  $a$

$T(s, a, s')$ : probability of transition from state  $s$  to  $s'$  when action  $a$  is taken

$R(s, a)$ : reward, takes in a (state, action) pair & returns a reward

$\gamma \in [0, 1]$ : discount factor, a scalar


$\pi(s)$ : policy, takes in state & return action

\* Goal of MDP is to find good policy

in 6.390:

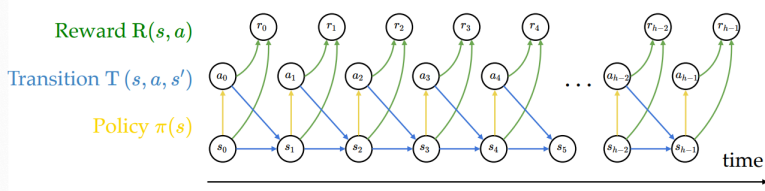
• S & A = small discrete sets

•  $s'$  &  $a'$  short-hand for next timestep state & action



Running example: Mario in a grid-world

- 9 possible states  $s$
- 4 possible actions  $a$ : {Up ↑, Down ↓, Left ←, Right →}
- (state, action) results in a transition  $T$  into a next state:
  - Normally, we get to the "intended" state;
    - E.g., in state (7), action "↑" gets to state (4)
  - If an action would take Mario out of the grid world, stay put;
    - E.g., in state (9), "←" gets back to state (9)
  - In state (6), action "↑" leads to two possibilities:
    - 20% chance to (2)
    - 80% chance to (3).



your future only depends on your current & not your past  
 every step & action determines the next future step

a trajectory (also called an experience or rollout) of horizon  $h$

$$\tau = (s_0, a_0, r_0, s_1, a_1, r_1, \dots, s_{h-1}, a_{h-1}, r_{h-1})$$

initial state

all depends on  $\pi$

- $a_t = \pi(s_t)$
- $r_t = R(s_t, a_t)$
- $T(s, a, s')$

[POLICY EVAL]

VALUE FUNC: long-term, sum & calc. expectation

Value functions:

$$V_h^{\pi}(s) := \mathbb{E} [R(s_0, \pi(s_0)) + \gamma R(s_1, \pi(s_1)) + \gamma^2 R(s_2, \pi(s_2)) + \gamma^3 R(s_3, \pi(s_3)) + \dots + \gamma^{h-1} R(s_{h-1}, \pi(s_{h-1}))]$$

$$= \mathbb{E} \left[ \sum_{t=0}^{h-1} \gamma^t R(s_t, \pi(s_t)) \mid s_0 = s, \pi \right] \quad (\text{eq. 1})$$

- $V_h^{\pi}(s)$  : expected sum of discounted rewards starting in state  $s$  and follow  $\pi$  for  $h$  steps
- Value is long-term; reward is immediate (one-time)
- Convention:  $V_0^{\pi}(s) = 0$  for all states
- The value expectation in 6.390 is only w.r.t. the transition probabilities  $T(s, a, s')$



evaluate  $V_h^\pi(s)$  under the "always-up" policy

states and one special transition:

|       |       |   |
|-------|-------|---|
| 1     | 2     | 3 |
| 20% ↗ | 80% ↘ |   |
| 4     | 5     | 6 |
| 7     | 8     | 9 |

rewards

|   |   |   |
|---|---|---|
| 0 | 0 | 0 |
| 0 | 0 | 0 |
| 0 | 0 | 0 |

- $\pi(s) = \uparrow$ ,  $\forall s$
- $\gamma = 0.9$

$$V_h^\pi(s) = \mathbb{E} \left[ \sum_{t=0}^{h-1} \gamma^t R(s_t, \uparrow) \mid s_0 = s \right]$$

$$= \mathbb{E} \left[ \underbrace{R(s_0, \uparrow) + \gamma R(s_1, \uparrow) + \dots + \gamma^{h-1} R(s_{h-1}, \uparrow)}_{h \text{ terms}} \right]$$

horizon  $h = 0$ : no step left

$$V_0^\pi(s) = 0$$

|   |   |   |
|---|---|---|
| 0 | 0 | 0 |
| 0 | 0 | 0 |
| 0 | 0 | 0 |

horizon  $h = 1$ : receive the rewards

$$V_1^\pi(s) = R(s, \uparrow)$$

|   |   |     |
|---|---|-----|
| 0 | 0 | 1   |
| 0 | 0 | -10 |
| 0 | 0 | 0   |



horizon  $h = 2$

states and one special transition:

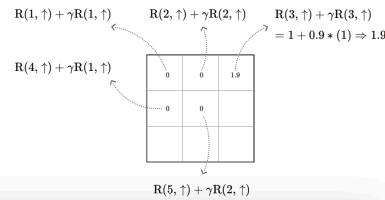
|       |       |   |
|-------|-------|---|
| 1     | 2     | 3 |
| 20% ↗ | 80% ↘ |   |
| 4     | 5     | 6 |
| 7     | 8     | 9 |

rewards

|   |   |   |
|---|---|---|
| 0 | 0 | 0 |
| 0 | 0 | 0 |
| 0 | 0 | 0 |

- $\pi(s) = \uparrow$ ,  $\forall s$
- $\gamma = 0.9$

$$V_2^\pi(s) = \mathbb{E} \left[ \underbrace{R(s_0, \uparrow) + \gamma R(s_1, \uparrow)}_{2 \text{ terms}} \right]$$



horizon  $h = 3$

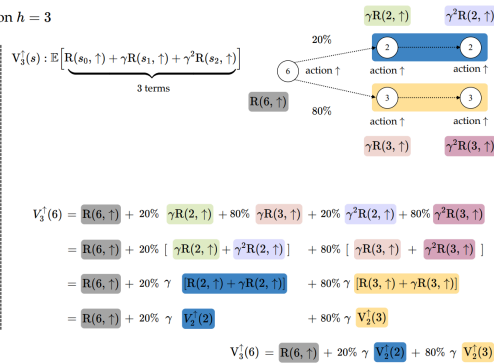
states and one special transition:

|       |       |   |
|-------|-------|---|
| 1     | 2     | 3 |
| 20% ↗ | 80% ↘ |   |
| 4     | 5     | 6 |
| 7     | 8     | 9 |

rewards

|   |   |   |
|---|---|---|
| 0 | 0 | 0 |
| 0 | 0 | 0 |
| 0 | 0 | 0 |

- $\pi(s) = \uparrow$ ,  $\forall s$
- $\gamma = 0.9$



horizon- $h$  value in state  $s$ : the expected sum of discounted rewards, starting in state  $s$  and following policy  $\pi$  for  $h$  steps.

$$(eq. 2) \quad V_h^\pi(s) = R(s, \pi(s)) + \gamma \sum_{s'} T(s, \pi(s), s') V_{h-1}^\pi(s')$$

the immediate reward for taking the policy-prescribed action  $\pi(s)$  in state  $s$ .

$(h - 1)$  horizon future value at a next state  $s'$

sum of future values weighted by the probability of reaching that next state  $s'$

discounted by  $\gamma$

group tail portions & recursively

Value functions converge as  $h \rightarrow \infty$

states and one special transition:

|       |       |   |
|-------|-------|---|
| 1     | 2     | 3 |
| 20% ↗ | 80% ↘ |   |
| 4     | 5     | 6 |
| 7     | 8     | 9 |

rewards

|   |   |   |
|---|---|---|
| 0 | 0 | 0 |
| 0 | 0 | 0 |
| 0 | 0 | 0 |

- $\pi(s) = \uparrow$ ,  $\forall s$
- $\gamma = 0.9$

| $V_{h0}^\pi(s)$ | $V_{h1}^\pi(s)$ | $V_{h\infty}^\pi(s)$ |
|-----------------|-----------------|----------------------|
| 0 0 9.98        | 0 0 9.98        | 0 0 10               |
| 0 0 -2.81       | 0 0 -2.81       | 0 0 -2.8             |
| 0 0 -2.53       | 0 0 -2.53       | 0 0 -2.52            |

- As we extend the horizon, value differences shrink
- because longer-term rewards are heavily discounted
- so, as  $h \rightarrow \infty$ , the value functions stop changing
- convergence can be seen, e.g., via  $V_{\infty}^\pi(3) = 1 + .9 + .9^2 + .9^3 + \dots = 10$

Typically,  $\gamma < 1$  to ensure  $V_{\infty}$  is finite.

die down to static vals

As horizon  $h \rightarrow \infty$ , the Bellman recursion becomes the Bellman equation

states and one special transition:

|       |       |   |
|-------|-------|---|
| 1     | 2     | 3 |
| 20% ↗ | 80% ↘ |   |
| 4     | 5     | 6 |
| 7     | 8     | 9 |

rewards

|   |   |   |
|---|---|---|
| 0 | 0 | 0 |
| 0 | 0 | 0 |
| 0 | 0 | 0 |

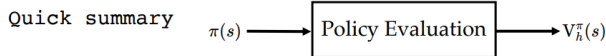
- $\pi(s) = \uparrow$ ,  $\forall s$
- $\gamma = 0.9$

$$\text{Recursion (finite } h) \quad (2) \quad V_h^\pi(s) = R(s, \pi(s)) + \gamma \sum_{s'} T(s, \pi(s), s') V_{h-1}^\pi(s')$$

$$\text{Equation (} h \rightarrow \infty) \quad (3) \quad V_{\infty}^\pi(s) = R(s, \pi(s)) + \gamma \sum_{s'} T(s, \pi(s), s') V_{\infty}^\pi(s')$$

A system of  $|\mathcal{S}|$  self-consistent linear equations, one for each state

| $V_{\infty}^\pi(s)$ | Equation  |
|---------------------|---|
| 0 0 10              | $V_{\infty}^\pi(3) = R(3, \uparrow) + \gamma V_{\infty}^\pi(3)$<br>$= 1 + .9 \times 10 \Rightarrow 10$  |
| 0 0 -2.8            | $V_{\infty}^\pi(6) = R(6, \uparrow) + \gamma [2 \times V_{\infty}^\pi(2) + 8 \times V_{\infty}^\pi(3)]$<br>$= -10 + .9 [2 \times 0 + 0.8 \times 10] \Rightarrow -2.8$ |
| 0 0 -2.52           |   |



Use the definition and sum up expected rewards:

1  $V_h^\pi(s) := \mathbb{E} \left[ \sum_{t=0}^{h-1} \gamma^t R(s_t, \pi(s_t)) \mid s_0 = s, \pi \right]$

Or, leverage the recursive structure:

2 finite-horizon Bellman recursions

$$V_h^\pi(s) = R(s, \pi(s)) + \gamma \sum_{s'} T(s, \pi(s), s') V_{h-1}^\pi(s')$$

3 infinite-horizon Bellman equations

$$V_\infty^\pi(s) = R(s, \pi(s)) + \gamma \sum_{s'} T(s, \pi(s), s') V_\infty^\pi(s')$$

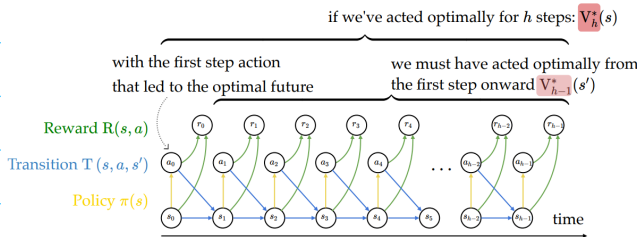
[POLICY OPTIMIZATION]

OPTIMAL POLICY  $\pi^*$ : policy that yields highest possible value  $V_h^*(s)$  from every state

• might not be unique

$$V_h^*(s) = \max_{\pi} V_h^\pi(s) = V_h^*(s), \forall s \in S$$

•  $V^*(s)$  defined over states, not actions



$$V_h^\pi(s) = R(s, \pi(s)) + \gamma \sum_{s'} T(s, \pi(s), s') V_{h-1}^\pi(s') \quad (\text{recall, eq. 2, for any policy})$$

$$V_h^*(s) = \max_a [R(s, a) + \gamma \sum_{s'} T(s, a, s') V_{h-1}^*(s')] \quad (\text{new, eq. 4, for an optimal policy})$$

$$V_h^*(s) = \max_a \underbrace{[R(s, a) + \gamma \sum_{s'} T(s, a, s') V_{h-1}^*(s')]}_{Q_h^*(s, a)} \quad (\text{eq. 4})$$

Define the optimal state-action value functions  $Q_h^*(s, a)$ :

the expected sum of discounted rewards, obtained by

- starting in state  $s$
- take action  $a$ , for one step
- act **optimally** thereafter for the remaining  $(h - 1)$  steps

$Q^*$  satisfies:

$$Q_h^*(s, a) = R(s, a) + \gamma \sum_{s'} T(s, a, s') \max_{a'} Q_{h-1}^*(s', a') \quad (\text{eq. 5})$$

$Q_h^*(s, a)$ : the value for

- starting in state  $s$ ,
- take action  $a$ , for one step
- act **optimally** thereafter for the remaining  $(h - 1)$  steps

States and one special transition:  $\gamma = 0.9$

|   |   |   |
|---|---|---|
| 1 | 2 | 3 |
| 4 | 5 | 6 |
| 7 | 8 | 9 |

Rewards:

|   |   |   |
|---|---|---|
| 0 | 0 | 1 |
| 0 | 0 | 1 |
| 0 | 0 | 1 |
| 0 | 0 | 1 |
| 0 | 0 | 1 |
| 0 | 0 | 1 |
| 0 | 0 | 1 |
| 0 | 0 | 1 |
| 0 | 0 | 1 |
| 0 | 0 | 1 |

Let's consider  $Q_2^*(6, \uparrow)$

- receive  $R(6, \uparrow)$
- act **optimally** at the next state  $s'$
- 20% chance,  $s' = 2$ , act optimally, get  $\max_{a'} Q_1^*(2, a')$
- 80% chance,  $s' = 3$ , act optimally, get  $\max_{a'} Q_1^*(3, a')$

$$Q_2^*(6, \uparrow) = R(6, \uparrow) + \gamma [0.2 \max_{a'} Q_1^*(2, a') + 0.8 \max_{a'} Q_1^*(3, a')]$$

$$= -10 + .9 [0.2 \times 0 + 0.8 \times 1] \Rightarrow -9.28$$

Value iteration: what we just did, iteratively invoke (eq. 5):

$$Q_h^*(s, a) = R(s, a) + \gamma \sum_{s'} T(s, a, s') \max_{a'} Q_{h-1}^*(s', a')$$

Value Iteration

1. for  $s \in \mathcal{S}, a \in \mathcal{A}$ :

2.  $Q_{\text{old}}(s, a) = 0$

3. while True:

4. for  $s \in \mathcal{S}, a \in \mathcal{A}$ :

5.  $Q_{\text{new}}(s, a) \leftarrow R(s, a) + \gamma \sum_{s'} T(s, a, s') \max_{a'} Q_{\text{old}}(s', a')$

6. if  $\max_{s, a} |Q_{\text{old}}(s, a) - Q_{\text{new}}(s, a)| < \epsilon$ :

7. return  $Q_{\text{new}}$

8.  $Q_{\text{old}} \leftarrow Q_{\text{new}}$

if run this block  $h$  times and break, then the returns are exactly  $Q_h^*$

$$Q_{\infty}^*(s, a)$$

$T(s, b, s')$  - action b matrix:

(cols = where you end up (s'))

|          |  |  |  |  |
|----------|--|--|--|--|
| from s=0 |  |  |  |  |
| from s=1 |  |  |  |  |
| from s=2 |  |  |  |  |
| from s=3 |  |  |  |  |

whole row: from state  $s$ , take action  $b$

Optimal policy easily extracted:  $\pi_h^*(s) = \arg \max_a Q_h^*(s, a)$

e.g. the best actions to take in state 5

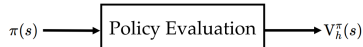
| $Q_1^*(s, a)$   | $Q_2^*(s, a)$ | $Q_3^*(s, a)$ | ... | $Q_{\infty}^*(s, a)$ |   |     |   |   |     |  |   |   |       |   |   |     |   |   |     |   |   |      |       |   |       |        |   |       |     |     |  |      |      |       |      |      |       |      |      |      |
|---|---------------|---------------|-----|----------------------|---|-----|---|---|-----|--|---|---|-------|---|---|-----|---|---|-----|---|---|------|-------|---|-------|--------|---|-------|-----|-----|--|------|------|-------|------|------|-------|------|------|------|
| <table border="1"><tr><td>0</td><td>0</td><td>1</td></tr><tr><td>0</td><td>0</td><td>1</td></tr><tr><td>0</td><td>0</td><td>1</td></tr></table>       | 0             | 0             | 1   | 0                    | 0 | 1   | 0 | 0 | 1   | <table border="1"><tr><td>0</td><td>0</td><td>1.9</td></tr><tr><td>0</td><td>0</td><td>1.9</td></tr><tr><td>0</td><td>0</td><td>-8</td></tr></table>   | 0 | 0 | 1.9   | 0 | 0 | 1.9 | 0 | 0 | -8  | <table border="1"><tr><td>0</td><td>0.81</td><td>2.71</td></tr><tr><td>0</td><td>1.71</td><td>2.71</td></tr><tr><td>0</td><td>-7.35</td><td></td></tr></table>    | 0 | 0.81 | 2.71  | 0 | 1.71  | 2.71   | 0 | -7.35 |     | ... | <table border="1"><tr><td>7.29</td><td>8.1</td><td>10</td></tr><tr><td>7.29</td><td>8.1</td><td>9</td></tr><tr><td>6.56</td><td>7.29</td><td>9.1</td></tr></table>         | 7.29 | 8.1  | 10    | 7.29 | 8.1  | 9     | 6.56 | 7.29 | 9.1  |
| 0   | 0             | 1             |     |                      |   |     |   |   |     |  |   |   |       |   |   |     |   |   |     |   |   |      |       |   |       |        |   |       |     |     |  |      |      |       |      |      |       |      |      |      |
| 0   | 0             | 1             |     |                      |   |     |   |   |     |  |   |   |       |   |   |     |   |   |     |   |   |      |       |   |       |        |   |       |     |     |  |      |      |       |      |      |       |      |      |      |
| 0   | 0             | 1             |     |                      |   |     |   |   |     |  |   |   |       |   |   |     |   |   |     |   |   |      |       |   |       |        |   |       |     |     |  |      |      |       |      |      |       |      |      |      |
| 0   | 0             | 1.9           |     |                      |   |     |   |   |     |  |   |   |       |   |   |     |   |   |     |   |   |      |       |   |       |        |   |       |     |     |  |      |      |       |      |      |       |      |      |      |
| 0   | 0             | 1.9           |     |                      |   |     |   |   |     |  |   |   |       |   |   |     |   |   |     |   |   |      |       |   |       |        |   |       |     |     |  |      |      |       |      |      |       |      |      |      |
| 0   | 0             | -8            |     |                      |   |     |   |   |     |  |   |   |       |   |   |     |   |   |     |   |   |      |       |   |       |        |   |       |     |     |  |      |      |       |      |      |       |      |      |      |
| 0   | 0.81          | 2.71          |     |                      |   |     |   |   |     |  |   |   |       |   |   |     |   |   |     |   |   |      |       |   |       |        |   |       |     |     |  |      |      |       |      |      |       |      |      |      |
| 0   | 1.71          | 2.71          |     |                      |   |     |   |   |     |  |   |   |       |   |   |     |   |   |     |   |   |      |       |   |       |        |   |       |     |     |  |      |      |       |      |      |       |      |      |      |
| 0   | -7.35         |               |     |                      |   |     |   |   |     |  |   |   |       |   |   |     |   |   |     |   |   |      |       |   |       |        |   |       |     |     |  |      |      |       |      |      |       |      |      |      |
| 7.29  | 8.1           | 10            |     |                      |   |     |   |   |     |  |   |   |       |   |   |     |   |   |     |   |   |      |       |   |       |        |   |       |     |     |  |      |      |       |      |      |       |      |      |      |
| 7.29  | 8.1           | 9             |     |                      |   |     |   |   |     |  |   |   |       |   |   |     |   |   |     |   |   |      |       |   |       |        |   |       |     |     |  |      |      |       |      |      |       |      |      |      |
| 6.56  | 7.29          | 9.1           |     |                      |   |     |   |   |     |  |   |   |       |   |   |     |   |   |     |   |   |      |       |   |       |        |   |       |     |     |  |      |      |       |      |      |       |      |      |      |
| <table border="1"><tr><td>0</td><td>0</td><td>-10</td></tr><tr><td>0</td><td>0</td><td>-10</td></tr><tr><td>0</td><td>0</td><td>-10</td></tr></table> | 0             | 0             | -10 | 0                    | 0 | -10 | 0 | 0 | -10 | <table border="1"><tr><td>0</td><td>0</td><td>-9.28</td></tr><tr><td>0</td><td>0</td><td>-9</td></tr><tr><td>0</td><td>0</td><td>-10</td></tr></table> | 0 | 0 | -9.28 | 0 | 0 | -9  | 0 | 0 | -10 | <table border="1"><tr><td>0</td><td>0.81</td><td>-8.47</td></tr><tr><td>0</td><td>-8.35</td><td>-18.35</td></tr><tr><td>0</td><td>0</td><td>-10</td></tr></table> | 0 | 0.81 | -8.47 | 0 | -8.35 | -18.35 | 0 | 0     | -10 | ... | <table border="1"><tr><td>7.29</td><td>8.1</td><td>-1.18</td></tr><tr><td>6.56</td><td>7.29</td><td>-1.07</td></tr><tr><td>5.9</td><td>6.56</td><td>-4.1</td></tr></table> | 7.29 | 8.1  | -1.18 | 6.56 | 7.29 | -1.07 | 5.9  | 6.56 | -4.1 |
| 0   | 0             | -10           |     |                      |   |     |   |   |     |  |   |   |       |   |   |     |   |   |     |   |   |      |       |   |       |        |   |       |     |     |  |      |      |       |      |      |       |      |      |      |
| 0   | 0             | -10           |     |                      |   |     |   |   |     |  |   |   |       |   |   |     |   |   |     |   |   |      |       |   |       |        |   |       |     |     |  |      |      |       |      |      |       |      |      |      |
| 0   | 0             | -10           |     |                      |   |     |   |   |     |  |   |   |       |   |   |     |   |   |     |   |   |      |       |   |       |        |   |       |     |     |  |      |      |       |      |      |       |      |      |      |
| 0   | 0             | -9.28         |     |                      |   |     |   |   |     |  |   |   |       |   |   |     |   |   |     |   |   |      |       |   |       |        |   |       |     |     |  |      |      |       |      |      |       |      |      |      |
| 0   | 0             | -9            |     |                      |   |     |   |   |     |  |   |   |       |   |   |     |   |   |     |   |   |      |       |   |       |        |   |       |     |     |  |      |      |       |      |      |       |      |      |      |
| 0   | 0             | -10           |     |                      |   |     |   |   |     |  |   |   |       |   |   |     |   |   |     |   |   |      |       |   |       |        |   |       |     |     |  |      |      |       |      |      |       |      |      |      |
| 0   | 0.81          | -8.47         |     |                      |   |     |   |   |     |  |   |   |       |   |   |     |   |   |     |   |   |      |       |   |       |        |   |       |     |     |  |      |      |       |      |      |       |      |      |      |
| 0   | -8.35         | -18.35        |     |                      |   |     |   |   |     |  |   |   |       |   |   |     |   |   |     |   |   |      |       |   |       |        |   |       |     |     |  |      |      |       |      |      |       |      |      |      |
| 0   | 0             | -10           |     |                      |   |     |   |   |     |  |   |   |       |   |   |     |   |   |     |   |   |      |       |   |       |        |   |       |     |     |  |      |      |       |      |      |       |      |      |      |
| 7.29  | 8.1           | -1.18         |     |                      |   |     |   |   |     |  |   |   |       |   |   |     |   |   |     |   |   |      |       |   |       |        |   |       |     |     |  |      |      |       |      |      |       |      |      |      |
| 6.56  | 7.29          | -1.07         |     |                      |   |     |   |   |     |  |   |   |       |   |   |     |   |   |     |   |   |      |       |   |       |        |   |       |     |     |  |      |      |       |      |      |       |      |      |      |
| 5.9   | 6.56          | -4.1          |     |                      |   |     |   |   |     |  |   |   |       |   |   |     |   |   |     |   |   |      |       |   |       |        |   |       |     |     |  |      |      |       |      |      |       |      |      |      |
| <table border="1"><tr><td>0</td><td>0</td><td>0</td></tr><tr><td>0</td><td>0</td><td>0</td></tr><tr><td>0</td><td>0</td><td>0</td></tr></table>       | 0             | 0             | 0   | 0                    | 0 | 0   | 0 | 0 | 0   | <table border="1"><tr><td>0</td><td>0</td><td>0</td></tr><tr><td>0</td><td>0</td><td>0</td></tr><tr><td>0</td><td>0</td><td>0</td></tr></table>        | 0 | 0 | 0     | 0 | 0 | 0   | 0 | 0 | 0   | <table border="1"><tr><td>0</td><td>0</td><td>0</td></tr><tr><td>0</td><td>0</td><td>0</td></tr><tr><td>0</td><td>0</td><td>0</td></tr></table>                   | 0 | 0    | 0     | 0 | 0     | 0      | 0 | 0     | 0   | ... | <table border="1"><tr><td>6.56</td><td>7.29</td><td>-1.07</td></tr><tr><td>5.9</td><td>6.56</td><td>5.9</td></tr><tr><td>5.9</td><td>6.56</td><td>5.9</td></tr></table>    | 6.56 | 7.29 | -1.07 | 5.9  | 6.56 | 5.9   | 5.9  | 6.56 | 5.9  |
| 0   | 0             | 0             |     |                      |   |     |   |   |     |  |   |   |       |   |   |     |   |   |     |   |   |      |       |   |       |        |   |       |     |     |  |      |      |       |      |      |       |      |      |      |
| 0   | 0             | 0             |     |                      |   |     |   |   |     |  |   |   |       |   |   |     |   |   |     |   |   |      |       |   |       |        |   |       |     |     |  |      |      |       |      |      |       |      |      |      |
| 0   | 0             | 0             |     |                      |   |     |   |   |     |  |   |   |       |   |   |     |   |   |     |   |   |      |       |   |       |        |   |       |     |     |  |      |      |       |      |      |       |      |      |      |
| 0   | 0             | 0             |     |                      |   |     |   |   |     |  |   |   |       |   |   |     |   |   |     |   |   |      |       |   |       |        |   |       |     |     |  |      |      |       |      |      |       |      |      |      |
| 0   | 0             | 0             |     |                      |   |     |   |   |     |  |   |   |       |   |   |     |   |   |     |   |   |      |       |   |       |        |   |       |     |     |  |      |      |       |      |      |       |      |      |      |
| 0   | 0             | 0             |     |                      |   |     |   |   |     |  |   |   |       |   |   |     |   |   |     |   |   |      |       |   |       |        |   |       |     |     |  |      |      |       |      |      |       |      |      |      |
| 0   | 0             | 0             |     |                      |   |     |   |   |     |  |   |   |       |   |   |     |   |   |     |   |   |      |       |   |       |        |   |       |     |     |  |      |      |       |      |      |       |      |      |      |
| 0   | 0             | 0             |     |                      |   |     |   |   |     |  |   |   |       |   |   |     |   |   |     |   |   |      |       |   |       |        |   |       |     |     |  |      |      |       |      |      |       |      |      |      |
| 0   | 0             | 0             |     |                      |   |     |   |   |     |  |   |   |       |   |   |     |   |   |     |   |   |      |       |   |       |        |   |       |     |     |  |      |      |       |      |      |       |      |      |      |
| 6.56  | 7.29          | -1.07         |     |                      |   |     |   |   |     |  |   |   |       |   |   |     |   |   |     |   |   |      |       |   |       |        |   |       |     |     |  |      |      |       |      |      |       |      |      |      |
| 5.9   | 6.56          | 5.9           |     |                      |   |     |   |   |     |  |   |   |       |   |   |     |   |   |     |   |   |      |       |   |       |        |   |       |     |     |  |      |      |       |      |      |       |      |      |      |
| 5.9   | 6.56          | 5.9           |     |                      |   |     |   |   |     |  |   |   |       |   |   |     |   |   |     |   |   |      |       |   |       |        |   |       |     |     |  |      |      |       |      |      |       |      |      |      |

- For finite  $h$ , optimal policy  $\pi_h^*$  depends on how many time steps are left
- When  $h \rightarrow \infty$ , time no longer matters, i.e., there exists a stationary  $\pi^*$

## Summary

- A Markov decision process  $(\mathcal{S}, \mathcal{A}, T, R, \gamma)$  is the mathematical framework for sequential decision-making and the foundation of reinforcement learning.
- To evaluate a given policy  $\pi$ , we compute state value functions  $V^\pi(s)$  via the Bellman recursion (finite horizon) or the Bellman equation (infinite horizon).
- To find an optimal policy, we compute  $Q^*(s, a)$  via the value iteration algorithm, then act greedily:  $\pi^*(s) = \arg \max_a Q^*(s, a)$ .

## 4127 REINFORCEMENT LEARNING



Use the definition and sum up expected rewards:

1  $V_h^\pi(s) := \mathbb{E} \left[ \sum_{t=0}^{h-1} \gamma^t R(s_t, \pi(s_t)) \mid s_0 = s, \pi \right]$

Or, leverage the recursive structure:

2 finite-horizon Bellman recursions

$$V_h^\pi(s) = R(s, \pi(s)) + \gamma \sum_{s'} T(s, \pi(s), s') V_{h-1}^\pi(s')$$

3 infinite-horizon Bellman equations

$$V_\infty^\pi(s) = R(s, \pi(s)) + \gamma \sum_{s'} T(s, \pi(s), s') V_\infty^\pi(s')$$

Recall

horizon- $h$  value in state  $s$ : the expected sum of discounted rewards, starting in state  $s$  and following policy  $\pi$  for  $h$  steps.

2  $V_h^\pi(s) = R(s, \pi(s)) + \gamma \sum_{s'} T(s, \pi(s), s') V_{h-1}^\pi(s')$

the immediate reward for taking the policy-prescribed action  $\pi(s)$  in state  $s$ .

$(h-1)$  horizon future value at a next state  $s'$

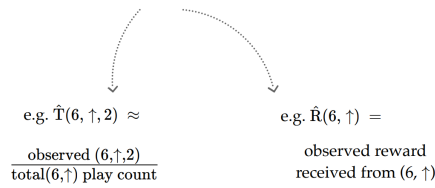
sum of future values weighted by the probability of reaching that next state  $s'$

discounted by  $\gamma$

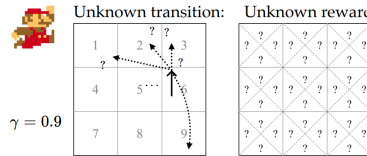
# Model-based RL: Learn the MDP tuple

1. Collect experiences  $(s, a, r, s')$

2. Estimate  $\hat{T}, \hat{R}$



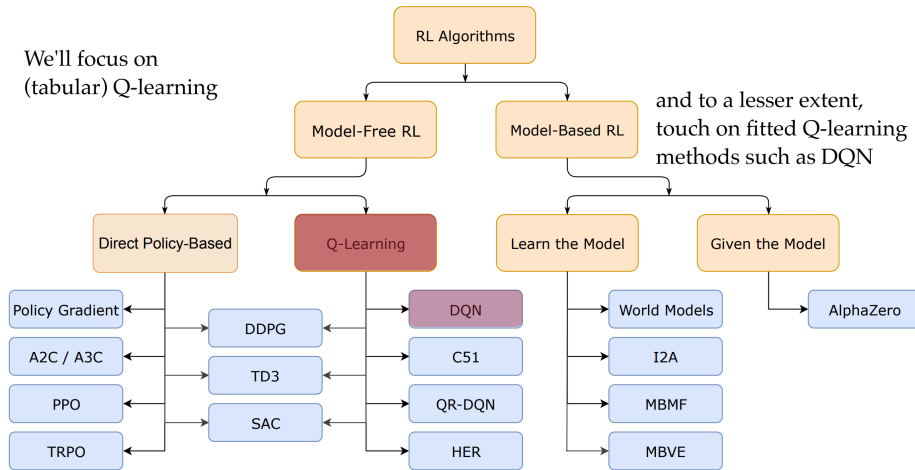
3. Solve  $\langle \mathcal{S}, \mathcal{A}, \hat{T}, \hat{R}, \gamma \rangle$  via e.g. *Value Iteration*



$\gamma = 0.9$

| $(s, a)$           | $r$ | $s'$ |
|--------------------|-----|------|
| (1, $\uparrow$ )   | 0   | 1    |
| (1, $\downarrow$ ) | 0   | 4    |
| ...                | ... | ...  |
| (3, $\uparrow$ )   | 1   | 3    |
| (3, $\downarrow$ ) | 1   | 6    |
| ...                | ... | ...  |
| (6, $\uparrow$ )   | -10 | 3    |
| (6, $\uparrow$ )   | -10 | 2    |
| ...                | ... | ...  |

We'll focus on (tabular) Q-learning



[A non-exhaustive, but useful taxonomy of algorithms in modern RL. Source] 14

## TABULAR Q-LEARNING:

once we have Q values, we basically have everything we need to get optimal policy

- Indeed, *value iteration* relies on having full access to R and T



$$Q_{\text{new}}(s, a) \leftarrow R(s, a) + \gamma \sum_{s'} T(s, a, s') \max_{a'} Q_{\text{old}}(s', a')$$

immediate reward      future value, starting in state  $s'$  and acting *optimally* for  $(h-1)$  steps

expected future value, weighted by the chance of landing in that particular future state  $s'$

- Without R and T, how about: execute  $(s, a)$ , observe  $r$  and  $s'$ , and update:

$$Q_{\text{new}}(s, a) \leftarrow r + \gamma \max_{a'} Q_{\text{old}}(s', a')$$

target

$$Q_{\text{new}}(s, a) \leftarrow r + \gamma \max_{a'} Q_{\text{old}}(s', a')$$

target

Whenever observe  $(6, \uparrow), -10, 3$ :

$$Q_{\text{new}}(6, \uparrow) \leftarrow -10 + \gamma \max_{a'} Q_{\text{old}}(3, a')$$

$$= -10 + 0.9 * 1 = -9.1$$

Whenever observe  $(6, \uparrow), -10, 2$ :

$$Q_{\text{new}}(6, \uparrow) \leftarrow -10 + \gamma \max_{a'} Q_{\text{old}}(2, a')$$

$$= -10 + 0.9 * 0 = -10$$



(we'll see why this fails)

Simply committing to the new target keeps "washing away" the old belief

$$Q_{\text{new}}(s, a) \leftarrow (1 - \alpha) Q_{\text{old}}(s, a) + \alpha \left( r + \gamma \max_{a'} Q_{\text{old}}(s', a') \right)$$

← downweigh w/  $\alpha$  and add in old belief

- Core update rule of Q-learning
- $\alpha \in [0, 1]$  : hyperparameter, trades off how much we trust new evidence versus old belief
- What we saw was  $\alpha = 1$ : trust new experience entirely
- Let's see an example using  $\alpha = 0.7$

$$Q_{\text{new}}(s, a) \leftarrow (1 - \alpha) Q_{\text{old}}(s, a) + \alpha (r + \gamma \max_{a'} Q_{\text{old}}(s', a'))$$

- To update a particular  $Q(s, a)$ , we need experience from actually taking that  $a$  in that  $s$ .
- If we never try a move, its Q-value never changes.
- So if we have one play left, which  $(s, a)$  should we try?

Depends on

★ [

- are we trying to "win" the game, or
- are we trying to get more accurate Q estimates

This is the fundamental dilemma in reinforcement learning... and in life!

Whether to *exploit* what we've already learned or *explore* to discover something better.

**ε-GREEDY ACTION SELECTION STRATEGY:**

- with probability  $\epsilon$ , choose action  $a \in \mathcal{A}$  uniformly @ random
- during learning, esp. in early stages, explore and observe diverse  $(s, a)$  consequences
- with probability  $(1 - \epsilon)$ , choose  $\text{argmax}_a Q_{\text{old}}(s, a)$
- during later stages, we can act more greedily w/ respect to current estimated Q values
- \* in beginning:  $\epsilon$  bigger (more exploratory), but smaller later (more greedy)
- \*  $\epsilon$  trades off exploration vs. exploitation

```

Value Iteration( $\mathcal{S}, \mathcal{A}, T, R, \gamma, \epsilon$ )
1. for  $s \in \mathcal{S}, a \in \mathcal{A}$ :
2.    $Q_{\text{old}}(s, a) = 0$ 
3. while True:
4.   for  $s \in \mathcal{S}, a \in \mathcal{A}$ :
5.      $Q_{\text{new}}(s, a) \leftarrow R(s, a) + \gamma \sum_{s'} T(s, a, s') \max_{a'} Q_{\text{old}}(s', a')$ 
6.   if  $\max_{s, a} |Q_{\text{old}}(s, a) - Q_{\text{new}}(s, a)| < \epsilon$ :
7.     return  $Q_{\text{new}}$ 
8.    $Q_{\text{old}} \leftarrow Q_{\text{new}}$ 
    
```

"calculating"

```

Q-Learning( $\mathcal{S}, \mathcal{A}, \gamma, \alpha, s_0, \text{max-iter}$ )
1.  $i = 0$ 
2. for  $s \in \mathcal{S}, a \in \mathcal{A}$ :
3.    $Q_{\text{old}}(s, a) = 0$ 
4.  $s \leftarrow s_0$ 
5. while  $i < \text{max-iter}$ :
6.    $a \leftarrow \text{select\_action}(s, Q_{\text{old}}(s, a))$ 
7.    $r, s' \leftarrow \text{execute}(a)$ 
8.    $Q_{\text{new}}(s, a) \leftarrow (1 - \alpha) Q_{\text{old}}(s, a) + \alpha (r + \gamma \max_{a'} Q_{\text{old}}(s', a'))$ 
9.    $s \leftarrow s'$ 
10.   $i \leftarrow (i + 1)$ 
11.   $Q_{\text{old}} \leftarrow Q_{\text{new}}$ 
12. return  $Q_{\text{new}}$ 
    
```

"learning" (estimating)

## Q-Learning ( $\mathcal{S}, \mathcal{A}, \gamma, \alpha, s_0, \text{max-iter}$ )

1.  $i = 0$
2. for  $s \in \mathcal{S}, a \in \mathcal{A}$ :
3.  $Q_{\text{old}}(s, a) = 0$
4.  $s \leftarrow s_0$
5. while  $i < \text{max-iter}$ :
6.  $a \leftarrow \text{select\_action}(s, Q_{\text{old}}(s, a))$
7.  $r, s' \leftarrow \text{execute}(a)$
8.  $Q_{\text{new}}(s, a) \leftarrow (1 - \alpha)Q_{\text{old}}(s, a) + \alpha(r + \gamma \max_{a'} Q_{\text{old}}(s', a'))$
9.  $s \leftarrow s'$
10.  $i \leftarrow (i + 1)$
11.  $Q_{\text{old}} \leftarrow Q_{\text{new}}$
12. return  $Q_{\text{new}}$

"learning"

- Remarkably,  $\rightarrow$  can converge to the true  $Q_{\infty}^*$
- Once converged, act greedily w.r.t  $Q^*$  again.
- But convergence can be extremely slow, and often not practical.
- Three hyperparameters:
  - discount factor  $\gamma$
  - learning rate  $\alpha$
  - exploration rate  $\epsilon$

all between 0 and 1

each controls some trade-off

<sup>1</sup> given that we visit all  $s, a$  infinitely often, and satisfy a decaying condition on the learning rate  $\alpha$ .

## FITTED Q-LEARNING:

• Q learning for large/continuous  $\mathcal{S}$  &  $\mathcal{A}$

• we can't keep a table, must approx.  $Q(s, a)$  w/ function  $Q_{\theta}(s, a)$

- Notice that the core update rule of Q-learning:

$$Q_{\text{new}}(s, a) \leftarrow (1 - \alpha)Q_{\text{old}}(s, a) + \alpha(r + \gamma \max_{a'} Q_{\text{old}}(s', a'))$$

is equivalent to:

$$Q_{\text{new}}(s, a) \leftarrow Q_{\text{old}}(s, a) + \alpha(r + \gamma \max_{a'} Q_{\text{old}}(s', a') - Q_{\text{old}}(s, a))$$

- Does this

$$\text{new belief} \leftarrow \text{old belief} + \text{learning rate} (\text{target} - \text{old belief})$$

remind you of something?

- Yes! Recall:

$$\theta_{\text{new}} \leftarrow \theta_{\text{old}} + \eta(\text{target} - \text{guess}_{\theta}) \nabla_{\theta} \text{guess}$$

Gradient update rule when minimizing  $(\text{target} - \text{guess}_{\theta})^2$

Q-learning is basically some version of gradient descent

- Same idea as before — we adjust our belief toward a target.

$$Q_{\text{new}}(s, a) \leftarrow Q_{\text{old}}(s, a) + \alpha(r + \gamma \max_{a'} Q_{\text{old}}(s', a') - Q_{\text{old}}(s, a))$$

$$\text{new belief} \leftarrow \text{old belief} + \text{learning rate} (\text{target} - \text{old belief})$$

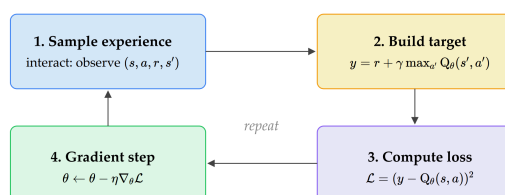
- Now Q is a function (e.g. a neural network) and  $\theta$  are its weights.

- Generalize tabular Q-learning for continuous state/action space:

1. parameterize  $Q_{\theta}(s, a)$
2. execute  $(s, a)$ , observe  $(r, s')$ , construct the target  $r + \gamma \max_{a'} Q_{\theta}(s', a')$
3. regress  $Q_{\theta}(s, a)$  against the target, i.e. update  $\theta$  via gradient-descent to minimize

$$(\text{target} - Q_{\theta}(s, a))^2$$

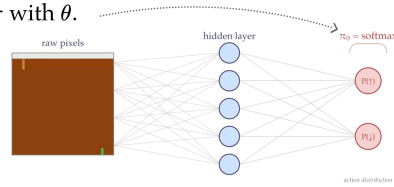
Fitted Q-learning: the training loop



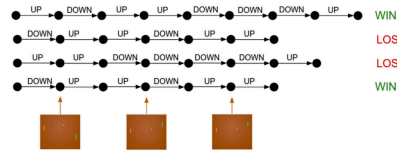
Same shape as supervised learning; except the "label"  $y$  is built from our own current estimate. (Bootstrapping.)

Policy gradients:

1. Parametrize the policy  $\pi$  with  $\theta$ .



2. Run  $\pi_\theta$  for a while. Keep track of the return.



3. Update  $\theta$  so "good" trajectories more likely to happen.

[image edited from: Andrej Karpathy]

Made rigorous by the *log-derivative trick* (REINFORCE). We will skip the derivation.

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim \pi_{\theta}} \left[ \sum_t \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \cdot R(\tau) \right]$$

Annotations in the diagram:  
 - "trajectory drawn under the current policy" points to  $\tau \sim \pi_{\theta}$   
 - "sum over timesteps" points to  $\sum_t$   
 - "push  $\theta$  to make the taken action more likely" points to  $\nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$   
 - "weight by how good the trajectory was" points to  $R(\tau)$

Sample trajectories, observe returns, nudge  $\theta$  to make actions in *good* trajectories more likely.

### How Much Information is the Machine Given during Learning?

- ▶ **"Pure" Reinforcement Learning (cherry)**
  - ▶ The machine predicts a scalar reward given once in a while.
  - ▶ **A few bits for some samples**
- ▶ **Supervised Learning (icing)**
  - ▶ The machine predicts a category or a few numbers for each input
  - ▶ Predicting human-supplied data
  - ▶ **10→10,000 bits per sample**
- ▶ **Self-Supervised Learning (cake génoise)**
  - ▶ The machine predicts any part of its input for any observed part.
  - ▶ Predicts future frames in videos
  - ▶ **Millions of bits per sample**



### Summary

- Last week: find good policies in a *known* MDP (high cumulative expected reward).
- In RL, the transitions and rewards are *unknown*. We still want a good policy; we get one by estimating Q.
- Exploration vs. exploitation: choose actions to gain reward while still learning.
- Q-learning converges to Q\* (with enough exploration).
- Deep Q-learning extends to large or continuous state spaces by parameterizing Q.

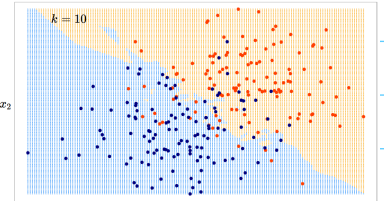
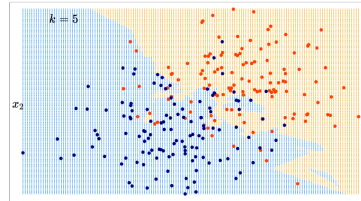
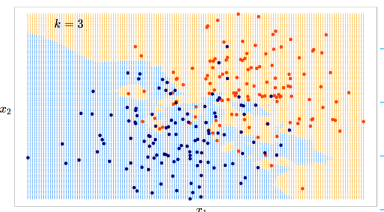
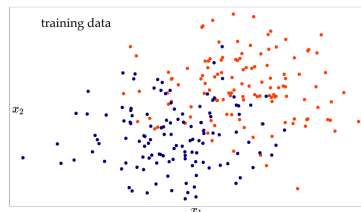
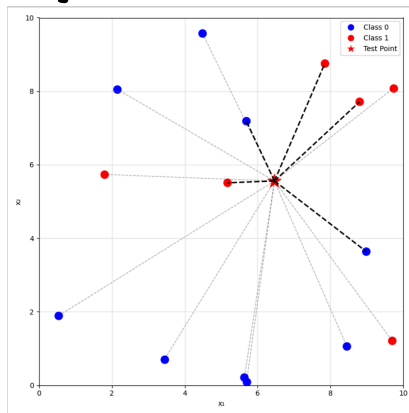
lecture 12 **NON-PARAMETRIC MODELS:** model doesn't assume a fixed functional form for  $h$ ; complexity grows with data

- interpretability (intuitive)
- insight
- speed
- adaptivity

**SIMILARITY-BASED METHODS:**

**k-nearest neighbor:**

prediction on new data using  $k = 5$  nearest neighbor



**be CAUTIONS when:**

- large  $n$  or  $d$  (curse of dimensionality)
- many irrelevant/noisy features
- features on diff. scales

Training: none. Just memorize the training set.

Predicting: for a new  $x_{new}$ , take the majority (class) or mean (regression) label of its  $k$  nearest neighbors.

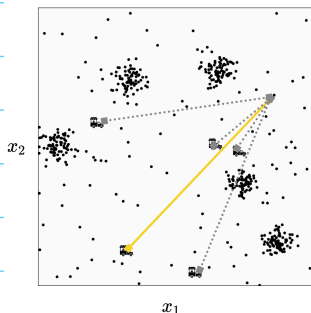
Parameters learned: the entire training set (its features and labels).  
becomes params

**Hyperparameters:**

- $k$ : number of neighbors considered.
- distance metric (typically Euclidean or Manhattan).
- tie-breaking scheme (typically at random).

**k means clustering: structural discovery w/o labels**

Per-person loss



Cost for person  $i$  to walk to truck  $j$ :

$$\|x^{(i)} - \mu^{(j)}\|^2$$

Let  $y^{(i)}$  be the truck person  $i$  walks to.

$\mathbf{1}\{y^{(i)} = j\}$  is 1 if assigned, 0 otherwise.

Person  $i$ 's loss:

$$\sum_{j=1}^k \mathbf{1}\{y^{(i)} = j\} \|x^{(i)} - \mu^{(j)}\|^2$$

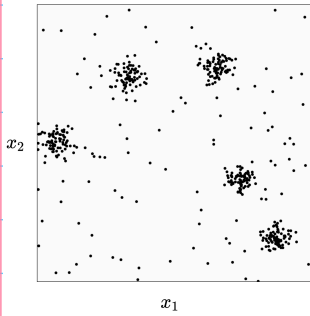
**k-means objective:**

cluster membership      centroid location

$$\sum_{i=1}^n \sum_{j=1}^k \mathbf{1}\{y^{(i)} = j\} \|x^{(i)} - \mu^{(j)}\|^2$$

sum over data points      sum over cluster

what we learn



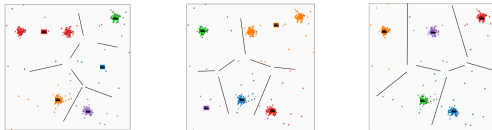
```

K-MEANS( $k, \tau, \{x^{(i)}\}_{i=1}^n$ )
1  $\mu, y =$  random initialization (rand. assign)
2 for  $t = 1$  to  $\tau$ 
3    $y_{old} = y$ 
4   for  $i = 1$  to  $n$ 
5      $y^{(i)} = \arg \min_j \|x^{(i)} - \mu^{(j)}\|^2$ 
6   for  $j = 1$  to  $k$ 
7      $\mu^{(j)} = \frac{1}{N_j} \sum_{i=1}^n \mathbf{1}(y^{(i)} = j) x^{(i)}$ 
8   if  $y == y_{old}$ 
9     break
10 return  $\mu, y$ 
    
```

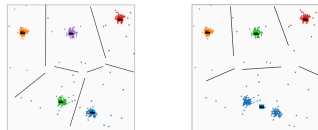
check until loc. chng  
cluster to/from

k-means is sensitive to initialization and to k

same data, three different initializations:



same data, two different choices of k:

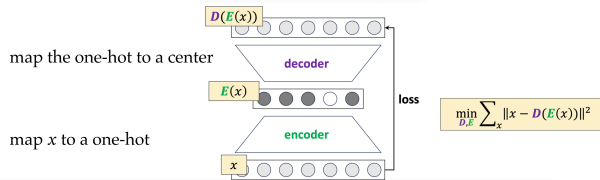


k-means interpreted as a discrete autoencoder

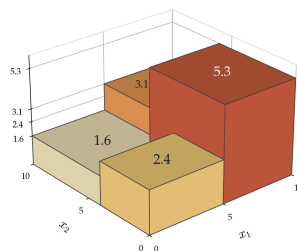
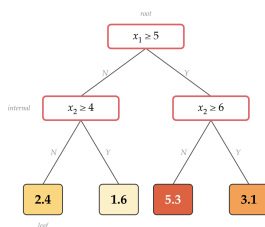


cluster center:

- not actual data points in MNIST
- learned representation and prototypes



Tree root node, internal node, leaf, and depth



Each internal node carries a split dimension and split value ( $x_j \geq \theta$ ).  
Each leaf corresponds to an axis-aligned region; the tree predicts the leaf's value there.

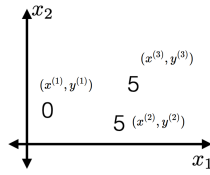
**tree-based method :**

BuildTree( $I, k, \mathcal{D}$ )

1. if  $|I| > k$
2. for each split dim  $j$  and split value  $s$
3. Set  $I_{j,s}^+ = \{i \in I \mid x_j^{(i)} \geq s\}$
4. Set  $I_{j,s}^- = \{i \in I \mid x_j^{(i)} < s\}$
5. Set  $\hat{y}_{j,s}^+ = \text{average}_{i \in I_{j,s}^+} y^{(i)}$
6. Set  $\hat{y}_{j,s}^- = \text{average}_{i \in I_{j,s}^-} y^{(i)}$
7. Set  $E_{j,s} = \sum_{i \in I_{j,s}^+} (y^{(i)} - \hat{y}_{j,s}^+)^2 + \sum_{i \in I_{j,s}^-} (y^{(i)} - \hat{y}_{j,s}^-)^2$
8. Set  $(j^*, s^*) = \arg \min_{j,s} E_{j,s}$
9. else
10. Set  $\hat{y} = \text{average}_{i \in I} y^{(i)}$
11. return Leaf(leave\_value= $\hat{y}$ )
12. return Node( $j^*, s^*, \text{BuildTree}(I_{j^*,s^*}^-, k), \text{BuildTree}(I_{j^*,s^*}^+, k)$ )

averaging pred.

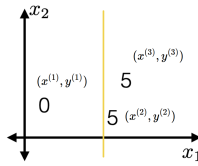
recursion



- Choose  $k = 2$
- BuildTree( $\{1, 2, 3\}; 2$ )
- Line 1 true

BuildTree( $I, k, \mathcal{D}$ )

1. if  $|I| > k$
2. for each split dim  $j$  and split value  $s$
3. Set  $I_{j,s}^+ = \{i \in I \mid x_j^{(i)} \geq s\}$
4. Set  $I_{j,s}^- = \{i \in I \mid x_j^{(i)} < s\}$
5. Set  $\hat{y}_{j,s}^+ = \text{average}_{i \in I_{j,s}^+} y^{(i)}$
6. Set  $\hat{y}_{j,s}^- = \text{average}_{i \in I_{j,s}^-} y^{(i)}$
7. Set  $E_{j,s} = \sum_{i \in I_{j,s}^+} (y^{(i)} - \hat{y}_{j,s}^+)^2 + \sum_{i \in I_{j,s}^-} (y^{(i)} - \hat{y}_{j,s}^-)^2$
8. Set  $(j^*, s^*) = \arg \min_{j,s} E_{j,s}$
9. else
10. Set  $\hat{y} = \text{average}_{i \in I} y^{(i)}$
11. return Leaf(leave\_value= $\hat{y}$ )
12. return Node( $j^*, s^*, \text{BuildTree}(I_{j^*,s^*}^-, k), \text{BuildTree}(I_{j^*,s^*}^+, k)$ )



- For this fixed  $(j, s)$ 
  - $I_{j,s}^+ = \{2, 3\}$
  - $I_{j,s}^- = \{1\}$
  - $\hat{y}_{j,s}^+ = 5$
  - $\hat{y}_{j,s}^- = 0$
  - $E_{j,s} = 0$

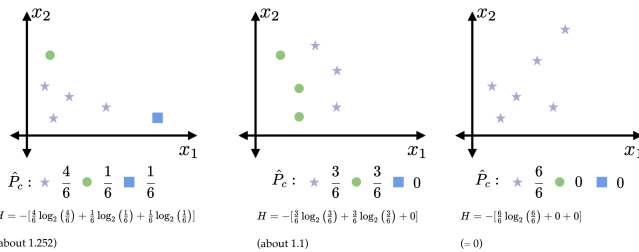
**CLASSIFICATION TREE:**

How to Grow a Classification Tree

1. Grow — start with all data at the root.  
Try candidate splits along each feature.
2. Split — pick the split that best separates the classes (lowest weighted entropy or most "homogeneous").
3. Recurse — treat each child region as a whole new dataset and run the algorithm again.
4. Stop — if a region is small or already pure (data within  $\leq$  leaf size, or most data share the same label), make it a leaf.

entropy  $H := - \sum_{c \in \text{class}} \hat{P}_c (\log_2 \hat{P}_c)$  empirical probability

for example:  $c$  iterates over 3 classes  $\star \bullet \blacksquare$



← Entropy is lower for more pure (less muddy) data

For classification

BuildTree( $I, k, \mathcal{D}$ )

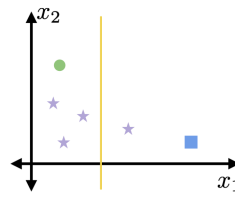
1. if  $|I| > k$
2. for each split dim  $j$  and split value  $s$
3. Set  $I_{j,s}^+ = \{i \in I \mid x_j^{(i)} \geq s\}$
4. Set  $I_{j,s}^- = \{i \in I \mid x_j^{(i)} < s\}$
5. Set  $\hat{y}_{j,s}^+ = \text{majority}_{i \in I_{j,s}^+} y^{(i)}$
6. Set  $\hat{y}_{j,s}^- = \text{majority}_{i \in I_{j,s}^-} y^{(i)}$
7. Set  $E_{j,s} = \frac{|I_{j,s}^+|}{|I|} \cdot H(I_{j,s}^+) + \frac{|I_{j,s}^-|}{|I|} \cdot H(I_{j,s}^-)$
8. Set  $(j^*, s^*) = \arg \min_{j,s} E_{j,s}$
9. else
10. Set  $\hat{y} = \text{majority}_{i \in I} y^{(i)}$
11. return Leaf(leave\_value= $\hat{y}$ )
12. return Node( $j^*, s^*, \text{BuildTree}(I_{j^*,s^*}^-, k), \text{BuildTree}(I_{j^*,s^*}^+, k)$ )

majority vote as regional prediction

weighted average entropy (WAE) as performance metric

**BuildTree( $I, k, \mathcal{D}$ )**

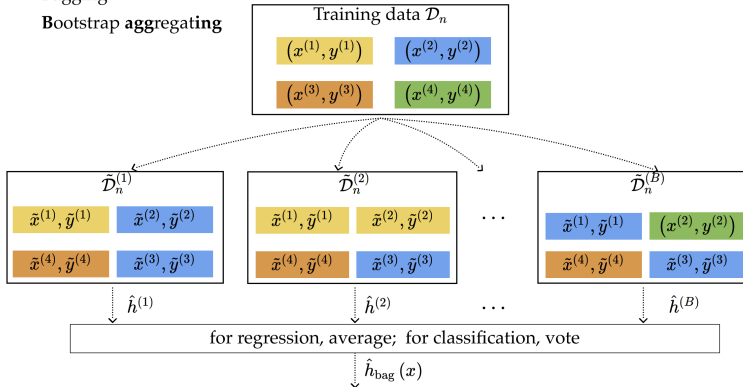
1. if  $|I| > k$
2. for each split dim  $j$  and split value  $s$
3. Set  $I_{j,s}^+ = \{i \in I \mid x_j^{(i)} \geq s\}$
4. Set  $I_{j,s}^- = \{i \in I \mid x_j^{(i)} < s\}$
5. Set  $\hat{y}_{j,s}^+ = \text{majority}_{i \in I_{j,s}^+} y^{(i)}$
6. Set  $\hat{y}_{j,s}^- = \text{majority}_{i \in I_{j,s}^-} y^{(i)}$
7. Set  $E_{j,s} = \frac{|I_{j,s}^-|}{|I|} \cdot H(I_{j,s}^-) + \frac{|I_{j,s}^+|}{|I|} \cdot H(I_{j,s}^+)$
8. Set  $(j^*, s^*) = \arg \min_{j,s} E_{j,s}$
9. else
10. Set  $\hat{y} = \text{majority}_{i \in I} y^{(i)}$
11. return Leaf(leave\_value= $\hat{y}$ )
12. return Node( $j^*, s^*, \text{BuildTree}(I_{j^*,s^*}^-, k), \text{BuildTree}(I_{j^*,s^*}^+, k)$ )



$$E_{j,s} = \frac{4}{6} \cdot H(I_{j,s}^-) + \frac{2}{6} \cdot H(I_{j,s}^+)$$

fraction of points to the left of the split      fraction of points to the right of the split

Bagging =  
Bootstrap aggregating



## Summary

- Non-parametric models let the data define structure, with few assumptions about functional form.
- $k$ -Nearest Neighbors: predict using nearby training points under a chosen distance metric.
- $k$ -Means: unsupervised grouping by distance, sensitive to initialization and the choice of  $k$ .
- Decision Trees: recursively split data into interpretable, flow-chart-like rules.
- Ensembles: combine many simple models so their votes become stronger, more stable predictors.