

Let's try to think about why PLA may work.

Let $n = n(t)$, according to the rule of PLA below, which formula is true?

$$\text{sign}(\mathbf{w}_t^T \mathbf{x}_n) \neq y_n, \quad \mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + y_n \mathbf{x}_n$$

- 1 $\mathbf{w}_{t+1}^T \mathbf{x}_n = y_n$
- 2 $\text{sign}(\mathbf{w}_{t+1}^T \mathbf{x}_n) = y_n$
- 3 $y_n \mathbf{w}_{t+1}^T \mathbf{x}_n \geq y_n \mathbf{w}_t^T \mathbf{x}_n$
- 4 $y_n \mathbf{w}_{t+1}^T \mathbf{x}_n < y_n \mathbf{w}_t^T \mathbf{x}_n$

Fun Time

Let's upper-bound T , the number of mistakes that PLA 'corrects'.

$$\text{Define } R^2 = \max_n \|\mathbf{x}_n\|^2 \quad \rho = \min_n y_n \frac{\mathbf{w}_f^T \mathbf{x}_n}{\|\mathbf{w}_f\|}$$

We want to show that $T \leq \square$. Express the upper bound \square by the two terms above.

- 1 R/ρ
- 2 R^2/ρ^2
- 3 R/ρ^2
- 4 ρ^2/R^2

Reference Answer: ②

The maximum value of $\frac{\mathbf{w}_f^T \mathbf{x}_n}{\|\mathbf{w}_f\| \|\mathbf{x}_n\|}$ is 1. Since T mistake corrections **increase the inner product by \sqrt{T} -constant**, the maximum number of corrected mistakes is $1/\text{constant}^2$.

Fun Time

Should we use pocket or PLA?

Since we do not know whether \mathcal{D} is linear separable in advance, we may decide to just go with pocket instead of PLA. If \mathcal{D} is actually linear separable, what's the difference between the two?

- 1 pocket on \mathcal{D} is slower than PLA
- 2 pocket on \mathcal{D} is faster than PLA
- 3 pocket on \mathcal{D} returns a better g in approximating f than PLA
- 4 pocket on \mathcal{D} returns a worse g in approximating f than PLA

Reference Answer: ①

Because pocket need to check whether \mathbf{w}_{t+1} is better than $\hat{\mathbf{w}}$ in each iteration, it is slower than PLA. On linear separable \mathcal{D} , $\mathbf{w}_{\text{POCKET}}$ is the same as \mathbf{w}_{PLA} , both making no mistakes.

Fun Time

This is a popular 'brain-storming' problem, with a claim that 2% of the world's cleverest population can crack its 'hidden pattern'.

$$(5, 3, 2) \rightarrow 151022, \quad (7, 2, 5) \rightarrow ?$$

It is like a 'learning problem' with $N = 1$, $\mathbf{x}_1 = (5, 3, 2)$, $y_1 = 151022$. Learn a hypothesis from the one example to predict on $\mathbf{x} = (7, 2, 5)$. What is your answer?

- ① 151026
- ② 143547
- ③ I need more examples to get the correct answer
- ④ there is no 'correct' answer

Reference Answer: ④

Following the same nature of the no-free-lunch problems discussed, we cannot hope to be correct under this 'adversarial' setting. BTW, ② is the designer's answer: the first two digits = $x_1 \cdot x_2$; the next two digits = $x_1 \cdot x_3$; the last two digits = $(x_1 \cdot x_2 + x_1 \cdot x_3 - x_2)$.

Fun Time

Let $\mu = 0.4$. Use Hoeffding's Inequality

$$\mathbb{P}[|\nu - \mu| > \epsilon] \leq 2 \exp(-2\epsilon^2 N)$$

to bound the probability that a sample of 10 marbles will have $\nu \leq 0.1$. What bound do you get?

- ① 0.67
- ② 0.40
- ③ 0.33
- ④ 0.05

Reference Answer: ③

Set $N = 10$ and $\epsilon = 0.3$ and you get the answer. BTW, ④ is the actual probability and Hoeffding gives only an upper bound to that.

Consider 4 hypotheses.

$$h_1(\mathbf{x}) = \text{sign}(x_1), \quad h_2(\mathbf{x}) = \text{sign}(x_2), \\ h_3(\mathbf{x}) = \text{sign}(-x_1), \quad h_4(\mathbf{x}) = \text{sign}(-x_2).$$

For any N and ϵ , which of the following statement is not true?

- ① the **BAD** data of h_1 and the **BAD** data of h_2 are exactly the same
- ② the **BAD** data of h_1 and the **BAD** data of h_3 are exactly the same
- ③ $\mathbb{P}_{\mathcal{D}}[\text{BAD for some } h_k] \leq 8 \exp(-2\epsilon^2 N)$
- ④ $\mathbb{P}_{\mathcal{D}}[\text{BAD for some } h_k] \leq 4 \exp(-2\epsilon^2 N)$

Reference Answer: ①

The important thing is to note that ② is true, which implies that ④ is true if you revisit the union bound. Similar ideas will be used to conquer the $M = \infty$ case.

Data size: how large do we need?

One way to use the inequality

$$\mathbb{P} [|E_{\text{in}}(g) - E_{\text{out}}(g)| > \epsilon] \leq \underbrace{2 \cdot M \cdot \exp(-2\epsilon^2 N)}_{\delta}$$

is to pick a tolerable difference ϵ as well as a tolerable **BAD** probability δ , and then gather data with size (N) large enough to achieve those tolerance criteria. Let $\epsilon = 0.1$, $\delta = 0.05$, and $M = 100$. What is the data size needed?

① 215

② 415

③ 615

④ 815

Reference Answer: ②

We can simply express N as a function of those 'known' variables. Then, the needed $N = \frac{1}{2\epsilon^2} \ln \frac{2M}{\delta}$.

What is the effective number of lines for five inputs $\in \mathbb{R}^2$?

① 14

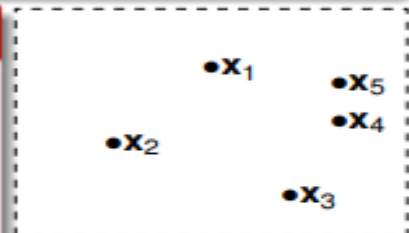
② 16

③ 22

④ 32

Reference Answer: ③

If you put the inputs roughly around a circle, you can then pick any consecutive inputs to be on one side of the line, and the other inputs to be on the other side. The procedure leads to effectively 22 kinds of lines, which is **much smaller than** $2^5 = 32$. You shall find it difficult to generate more kinds by varying the inputs, and we will give a formal proof in future lectures.



Fun Time

Consider positive **and negative** rays as \mathcal{H} , which is equivalent to the perceptron hypothesis set in 1D. The hypothesis set is often called '**decision stump**' to describe the shape of its hypotheses. What is the growth function $m_{\mathcal{H}}(N)$?

1 N

2 $N + 1$

3 $2N$

4 2^N

Reference Answer: 3

Two dichotomies when threshold in each of the $N - 1$ 'internal' spots; two dichotomies for the all- \circ and all- \times cases.

3

Fun Time

Consider positive **and negative** rays as \mathcal{H} , which is equivalent to the perceptron hypothesis set in 1D. As discussed in an earlier quiz question, the growth function $m_{\mathcal{H}}(N) = 2N$. What is the minimum break point for \mathcal{H} ?

1 1

2 2

3 3

4 4

Reference Answer: 3

At $k = 3$, $m_{\mathcal{H}}(k) = 6$ while $2^k = 8$.

+

Fun Time

When minimum break point $k = 1$, what is the maximum possible $m_{\mathcal{H}}(N)$ when $N = 3$?

- ① 1 ② 2 ③ 4 ④ 8

Reference Answer: ①

Because $k = 1$, the hypothesis set cannot even shatter one point. Thus, every 'column' of the table cannot contain both \circ and \times . Then, after including the first dichotomy, it is not possible to include any other different dichotomy. Thus, the maximum possible $m_{\mathcal{H}}(N)$ is 1.

x_1	x_2	x_3
\circ	\times	\circ
\circ	\times	\times

Fun Time

For the 2D perceptrons, which of the following claim is true?

- ① minimum break point $k = 2$
 ② $m_{\mathcal{H}}(4) = 15$
 ③ $m_{\mathcal{H}}(N) < B(N, k)$ when $N = k =$ minimum break point
 ④ $m_{\mathcal{H}}(N) > B(N, k)$ when $N = k =$ minimum break point

Reference Answer: ③

As discussed previously, minimum break point for 2D perceptrons is 2, with $m_{\mathcal{H}}(4) = 14$. Also, note that $B(4, 2) = 15$. So bounding function $B(N, k)$ can be 'loose' in bounding $m_{\mathcal{H}}(N)$.

Fun Time

For 1D perceptrons (positive and negative rays), we know that $m_{\mathcal{H}}(N) = 2N$. Let k be the minimum break point. Which of the following is not true?

- ① $k = 3$
- ② for some integers $N > 0$, $m_{\mathcal{H}}(N) = \sum_{i=0}^{k-1} \binom{N}{i}$
- ③ for all integers $N > 0$, $m_{\mathcal{H}}(N) = \sum_{i=0}^{k-1} \binom{N}{i}$
- ④ for all integers $N > 2$, $m_{\mathcal{H}}(N) < \sum_{i=0}^{k-1} \binom{N}{i}$

Reference Answer: ③

The proof is generally trivial by listing the definitions. For ②, $N = 1$ or 2 gives the equality. One thing to notice is ④: the upper bound can be 'loose'.

Fun Time

For positive rays, $m_{\mathcal{H}}(N) = N + 1$. Plug it into the VC bound for $\epsilon = 0.1$ and $N = 10000$. What is VC bound of BAD events?

$$\mathbb{P}[\exists h \in \mathcal{H} \text{ s.t. } |E_{\text{in}}(h) - E_{\text{out}}(h)| > \epsilon] \leq 4m_{\mathcal{H}}(2N) \exp\left(-\frac{1}{8}\epsilon^2 N\right)$$

- ① 2.77×10^{-87}
- ② 5.54×10^{-83}
- ③ 2.98×10^{-1}
- ④ 2.29×10^2

Reference Answer: ③

Simple calculation. Note that the BAD probability bound is not very small even with 10000 examples.

If there is a set of N inputs that cannot be shattered by \mathcal{H} . Based only on this information, what can we conclude about $d_{VC}(\mathcal{H})$?

- ① $d_{VC}(\mathcal{H}) > N$
- ② $d_{VC}(\mathcal{H}) = N$
- ③ $d_{VC}(\mathcal{H}) < N$
- ④ no conclusion can be made

Reference Answer: ④

It is possible that there is another set of N inputs that can be shattered, which means $d_{VC} \geq N$. It is also possible that no set of N input can be shattered, which means $d_{VC} < N$. Neither cases can be ruled out by one non-shattering set.

Extra Fun Time

What statement below shows that $d_{VC} \geq d + 1$?

- ① There are some $d + 1$ inputs we can shatter.
- ② We can shatter any set of $d + 1$ inputs.
- ③ There are some $d + 2$ inputs we cannot shatter.
- ④ We cannot shatter any set of $d + 2$ inputs.

Reference Answer: ①

d_{VC} is the maximum that $m_{\mathcal{H}}(N) = 2^N$, and $m_{\mathcal{H}}(N)$ is the most number of dichotomies of N inputs. So if we can find 2^{d+1} dichotomies on some $d + 1$ inputs, $m_{\mathcal{H}}(d + 1) = 2^{d+1}$ and hence $d_{VC} \geq d + 1$.

Extra Fun Time

What statement below shows that $d_{VC} \leq d + 1$?

- ① There are some $d + 1$ inputs we can shatter.
- ② We can shatter any set of $d + 1$ inputs.
- ③ There are some $d + 2$ inputs we cannot shatter.
- ④ We cannot shatter any set of $d + 2$ inputs.

Reference Answer: ④

d_{VC} is the maximum that $m_{\mathcal{H}}(N) = 2^N$, and $m_{\mathcal{H}}(N)$ is the most number of dichotomies of N inputs. So if we cannot find 2^{d+2} dichotomies on *any* $d + 2$ inputs (i.e. break point), $m_{\mathcal{H}}(d + 2) < 2^{d+2}$ and hence $d_{VC} < d + 2$. That is, $d_{VC} \leq d + 1$.

Fun Time

Based on the proof above, what is d_{VC} of 1126-D perceptrons?

- ① 1024
- ② 1126
- ③ 1127
- ④ 6211

Reference Answer: ③

Well, **too much fun for this section! :-)**

Fun Time

Origin-crossing Hyperplanes are essentially perceptrons with w_0 fixed at 0. Make a guess about the d_{VC} of origin-crossing hyperplanes in \mathbb{R}^d .

- ① 1
- ② d
- ③ $d + 1$
- ④ ∞

Reference Answer: ②

The proof is almost the same as proving the d_{VC} for usual perceptrons, but it is the **intuition** ($d_{VC} \approx \# \text{free parameters}$) that you shall use to answer this quiz.

Consider the VC Bound below. How can we decrease the probability of getting **BAD** data?

$$\mathbb{P}_{\mathcal{D}} \left[|E_{\text{in}}(g) - E_{\text{out}}(g)| > \epsilon \right] \leq 4(2N)^{d_{VC}} \exp \left(-\frac{1}{8} \epsilon^2 N \right)$$

- ① decrease model complexity d_{VC}
- ② increase data size N a lot
- ③ increase generalization error tolerance ϵ
- ④ all of the above

Reference Answer: ④

**Congratulations on being
Master of VC bound! :-)**

Fun Time

Let's revisit PLA/pocket. Which of the following claim is true?

- ① In practice, we should try to compute if \mathcal{D} is linear separable before deciding to use PLA.
- ② If we know that \mathcal{D} is not linear separable, then the target function f must not be a linear function.
- ③ If we know that \mathcal{D} is linear separable, then the target function f must be a linear function.
- ④ None of the above

Reference Answer: ④

① After computing if \mathcal{D} is linear separable, we shall know \mathbf{w}^* and then there is no need to use PLA. ② What about noise? ③ What about 'sampling luck'? :-)

Consider the following $P(y|\mathbf{x})$ and $\text{err}(\tilde{y}, y) = |\tilde{y} - y|$. Which of the following is the ideal mini-target $f(\mathbf{x})$?

$$P(y = 1|\mathbf{x}) = 0.10, P(y = 2|\mathbf{x}) = 0.35, \\ P(y = 3|\mathbf{x}) = 0.15, P(y = 4|\mathbf{x}) = 0.40.$$

- ① 2.5 = average within $\mathcal{Y} = \{1, 2, 3, 4\}$
- ② 2.85 = weighted mean from $P(y|\mathbf{x})$
- ③ 3 = weighted median from $P(y|\mathbf{x})$
- ④ 4 = $\text{argmax } P(y|\mathbf{x})$

Reference Answer: ③

For the 'absolute error', the weighted median provably results in the minimum average err.

Consider err below for CIA. What is $E_{in}(g)$ when using this err?

		g	
		+1	-1
f	+1	0	1
	-1	1000	0

① $\frac{1}{N} \sum_{n=1}^N \mathbb{I}[y_n \neq g(\mathbf{x}_n)]$
 ② $\frac{1}{N} \left(\sum_{y_n=+1} \mathbb{I}[y_n \neq g(\mathbf{x}_n)] + 1000 \sum_{y_n=-1} \mathbb{I}[y_n \neq g(\mathbf{x}_n)] \right)$
 ③ $\frac{1}{N} \left(\sum_{y_n=+1} \mathbb{I}[y_n \neq g(\mathbf{x}_n)] - 1000 \sum_{y_n=-1} \mathbb{I}[y_n \neq g(\mathbf{x}_n)] \right)$
 ④ $\frac{1}{N} \left(1000 \sum_{y_n=+1} \mathbb{I}[y_n \neq g(\mathbf{x}_n)] + \sum_{y_n=-1} \mathbb{I}[y_n \neq g(\mathbf{x}_n)] \right)$

Reference Answer: ②

When $y_n = -1$, the false positive made on such (\mathbf{x}_n, y_n) is penalized 1000 times more!

Consider the CIA cost matrix. If there are 10 examples with $y_n = -1$ (intruder) and 999,990 examples with $y_n = +1$ (you). What would $E_{in}^w(h)$ be for a constant $h(\mathbf{x})$ that always returns +1?

		$h(\mathbf{x})$	
		+1	-1
y	+1	0	1
	-1	1000	0

① 0.001
 ② 0.01
 ③ 0.1
 ④ 1

Reference Answer: ②

While the quiz is a simple evaluation, it is not uncommon that the data is very **unbalanced** for such an application. Properly 'setting' the weights can be used to avoid the lazy constant prediction.

Fun Time

Based on our discussion, for data of fixed size, which of the following situation is relatively of the lowest risk of overfitting?

- ① small noise, fitting from small d_{VC} to median d_{VC}
- ② small noise, fitting from small d_{VC} to large d_{VC}
- ③ large noise, fitting from small d_{VC} to median d_{VC}
- ④ large noise, fitting from small d_{VC} to large d_{VC}

Reference Answer: ①

Two causes of overfitting are noise and excessive d_{VC} . So if both are relatively 'under control', the risk of overfitting is smaller.

Fun Time

When having limited data, in which of the following case would learner R perform better than learner O ?

- ① limited data from a 10-th order target function with some noise
- ② limited data from a 1126-th order target function with no noise
- ③ limited data from a 1126-th order target function with some noise
- ④ all of the above

Reference Answer: ④

We discussed about ① and ②, but you shall be able to 'generalize' :-) that R also wins in the more difficult case of ③.

Fun Time

Consider the target function being $\sin(1126x)$ for $x \in [0, 2\pi]$. When x is uniformly sampled from the range, and we use all possible linear hypotheses $h(x) = w \cdot x$ to approximate the target function with respect to the squared error, what is the level of deterministic noise for each x ?

- ① $|\sin(1126x)|$
- ② $|\sin(1126x) - x|$
- ③ $|\sin(1126x) + x|$
- ④ $|\sin(1126x) - 1126x|$

Reference Answer: ①

You can try a few different w and convince yourself that the best hypothesis h^* is $h^*(x) = 0$. The deterministic noise is the difference between f and h^* .

Fun Time

Assume we know that $f(x)$ is symmetric for some 1D regression application. That is, $f(x) = f(-x)$. One possibility of using the knowledge is to consider symmetric hypotheses only. On the other hand, you can also generate virtual examples from the original data $\{(x_n, y_n)\}$ as hints. What virtual examples suit your needs best?

- ① $\{(x_n, -y_n)\}$
- ② $\{(-x_n, -y_n)\}$
- ③ $\{(-x_n, y_n)\}$
- ④ $\{(2x_n, 2y_n)\}$

Reference Answer: ③

We want the virtual examples to encode the invariance when $x \rightarrow -x$.

For $Q \geq 1$, which of the following hypothesis (weight vector $\mathbf{w} \in \mathbb{R}^{Q+1}$) is not in the regularized hypothesis set $\mathcal{H}(1)$?

- ① $\mathbf{w}^T = [0, 0, \dots, 0]$
- ② $\mathbf{w}^T = [1, 0, \dots, 0]$
- ③ $\mathbf{w}^T = [1, 1, \dots, 1]$
- ④ $\mathbf{w}^T = [\sqrt{\frac{1}{Q+1}}, \sqrt{\frac{1}{Q+1}}, \dots, \sqrt{\frac{1}{Q+1}}]$

Reference Answer: ③

The squared length of \mathbf{w} in ③ is $Q + 1$, which is not ≤ 1 .

Fun Time

When would \mathbf{w}_{REG} equal \mathbf{w}_{LIN} ?

- ① $\lambda = 0$
- ② $C = \infty$
- ③ $C \geq \|\mathbf{w}_{\text{LIN}}\|^2$
- ④ all of the above

Reference Answer: ④

① and ② shall be easy; ③ means that there are effectively no constraint on \mathbf{w} , hence the equivalence.

Fun Time

Consider the weight-decay regularization with regression. When increasing λ in \mathcal{A} , what would happen with $d_{\text{EFF}}(\mathcal{H}, \mathcal{A})$?

- ① $d_{\text{EFF}} \uparrow$
- ② $d_{\text{EFF}} \downarrow$
- ③ $d_{\text{EFF}} = d_{\text{VC}}(\mathcal{H})$ and does not depend on λ
- ④ $d_{\text{EFF}} = 1126$ and does not depend on λ

Reference Answer: ②

larger λ
 \iff smaller C
 \iff smaller $\mathcal{H}(C)$
 \iff smaller d_{EFF}

Fun Time

Consider using a regularizer $\Omega(\mathbf{w}) = \sum_{q=0}^Q 2^q w_q^2$ to work with Legendre polynomial regression. Which kind of hypothesis does the regularizer prefer?

- ① symmetric polynomials satisfying $h(x) = h(-x)$
- ② low-dimensional polynomials
- ③ high-dimensional polynomials
- ④ no specific preference

Reference Answer: ②

There is a higher 'penalty' for higher-order terms, and hence the regularizer prefers low-dimensional polynomials.

4

For $\mathcal{X} = \mathbb{R}^d$, consider two hypothesis sets, \mathcal{H}_+ and \mathcal{H}_- . The first hypothesis set contains all perceptrons with $w_1 \geq 0$, and the second hypothesis set contains all perceptrons with $w_1 \leq 0$. Denote g_+ and g_- as the minimum- E_{in} hypothesis in each hypothesis set, respectively. Which statement below is true?

- ① If $E_{\text{in}}(g_+) < E_{\text{in}}(g_-)$, then g_+ is the minimum- E_{in} hypothesis of all perceptrons in \mathbb{R}^d .
- ② If $E_{\text{test}}(g_+) < E_{\text{test}}(g_-)$, then g_+ is the minimum- E_{test} hypothesis of all perceptrons in \mathbb{R}^d .
- ③ The two hypothesis sets are disjoint.
- ④ None of the above

Reference Answer: ①

Note that the two hypothesis sets are not disjoint (sharing ' $w_1 = 0$ ' perceptrons) but their union is all perceptrons.

For a learning model that takes N^2 seconds of training when using N examples, what is the total amount of seconds needed when running the whole validation procedure with $K = \frac{N}{5}$ on 25 such models with different parameters to get the final g_{m^*} ?

- 1 $6N^2$
- 2 $17N^2$
- 3 $25N^2$
- 4 $26N^2$

Reference Answer: 2

To get all the g_m^- , we need $\frac{16}{25}N^2 \cdot 25$ seconds. Then to get g_{m^*} , we need another N^2 seconds. So in total we need $17N^2$ seconds.

Fun Time

Consider three examples (\mathbf{x}_1, y_1) , (\mathbf{x}_2, y_2) , (\mathbf{x}_3, y_3) with $y_1 = 1$, $y_2 = 5$, $y_3 = 7$. If we use E_{loocv} to estimate the performance of a learning algorithm that predicts with the average y value of the data set—the optimal constant prediction with respect to the squared error. What is E_{loocv} (squared error) of the algorithm?

- 1 0
- 2 $\frac{56}{9}$
- 3 $\frac{60}{9}$
- 4 14

Reference Answer: 4

This is based on a simple calculation of $e_1 = (1 - 6)^2$, $e_2 = (5 - 4)^2$, $e_3 = (7 - 3)^2$.

For a learning model that takes N^2 seconds of training when using N examples, what is the total amount of seconds needed when running 10-fold cross validation on 25 such models with different parameters to get the final g_{m^*} ?

- ① $\frac{47}{2} N^2$
- ② $47 N^2$
- ③ $\frac{407}{2} N^2$
- ④ $407 N^2$

Reference Answer: ③

To get all the E_{cv} , we need $\frac{81}{100} N^2 \cdot 10 \cdot 25$ seconds. Then to get g_{m^*} , we need another N^2 seconds. So in total we need $\frac{407}{2} N^2$ seconds.

Fun Time

For a deep NNet for written character recognition from raw pixels, which type of features are more likely extracted after the first hidden layer?

- ① pixels
- ② strokes
- ③ parts
- ④ digits

Reference Answer: ②

Simple strokes are likely the 'next-level' features that can be extracted from raw pixels.

Fun Time

Suppose training a d - \tilde{d} - d autoencoder with backprop takes approximately $c \cdot d \cdot \tilde{d}$ seconds. Then, what is the total number of seconds needed for pre-training a d - $d^{(1)}$ - $d^{(2)}$ - $d^{(3)}$ -1 deep NNet?

- ① $c (d + d^{(1)} + d^{(2)} + d^{(3)} + 1)$
- ② $c (d \cdot d^{(1)} \cdot d^{(2)} \cdot d^{(3)} \cdot 1)$
- ③ $c (d d^{(1)} + d^{(1)} d^{(2)} + d^{(2)} d^{(3)} + d^{(3)})$
- ④ $c (d d^{(1)} \cdot d^{(1)} d^{(2)} \cdot d^{(2)} d^{(3)} \cdot d^{(3)})$

Reference Answer: ③

Each $c \cdot d^{(\ell-1)} \cdot d^{(\ell)}$ represents the time for pre-training with one autoencoder to determine one layer of the weights.

Fun Time

Which of the following cannot be viewed as a regularization technique?

- ① hint the model with artificially-generated noisy data
- ② stop gradient descent early
- ③ add a weight elimination regularizer
- ④ all the above are regularization techniques

Reference Answer: ④

① is our new friend for regularization, while
② and ③ are old friends.

Fun Time

When solving the optimization problem

$$\max_{\mathbf{v}} \sum_{n=1}^N \mathbf{v}^T \mathbf{x}_n \mathbf{x}_n^T \mathbf{v} \text{ subject to } \mathbf{v}^T \mathbf{v} = 1,$$

we know that the optimal \mathbf{v} is the 'topmost' eigenvector that corresponds to the 'topmost' eigenvalue λ of $\mathbf{X}^T \mathbf{X}$. Then, what is the optimal objective value of the optimization problem?

- ① λ^1
- ② λ^2
- ③ λ^3
- ④ λ^4

Reference Answer: ①

The objective value of the optimization problem is simply $\mathbf{v}^T \mathbf{X}^T \mathbf{X} \mathbf{v}$, which is $\lambda \mathbf{v}^T \mathbf{v}$ and you know what $\mathbf{v}^T \mathbf{v}$ must be.

Fun Time

Let $g_0(\mathbf{x}) = +1$. Which of the following $(\alpha_0, \alpha_1, \alpha_2)$ allows

$G(\mathbf{x}) = \text{sign} \left(\sum_{t=0}^2 \alpha_t g_t(\mathbf{x}) \right)$ to implement $\text{OR}(g_1, g_2)$?

- ① $(-3, +1, +1)$
- ② $(-1, +1, +1)$
- ③ $(+1, +1, +1)$
- ④ $(+3, +1, +1)$

Reference Answer: ③

You can easily verify with all four possibilities of $(g_1(\mathbf{x}), g_2(\mathbf{x}))$.

Fun Time

How many weights $\{w_{ij}^{(l)}\}$ are there in a 3-5-1 NNet?

- ① 9
- ② 15
- ③ 20
- ④ 26

Reference Answer: ④

There are $(3 + 1) \times 5$ weights in $w_{ij}^{(1)}$, and $(5 + 1) \times 1$ weights in $w_{jk}^{(2)}$.

Fun Time

According to $\frac{\partial e_n}{\partial w_{i1}^{(L)}} = -2 \left(y_n - s_i^{(L)} \right) \cdot \left(x_i^{(L-1)} \right)$ when would $\frac{\partial e_n}{\partial w_{i1}^{(L)}} = 0$?

- ① $y_n = s_i^{(L)}$
- ② $x_i^{(L-1)} = 0$
- ③ $s_i^{(L-1)} = 0$
- ④ all of the above

Reference Answer: ④

Note that $x_i^{(L-1)} = \tanh(s_i^{(L-1)}) = 0$ if and only if $s_i^{(L-1)} = 0$.

Fun Time

For the weight elimination regularizer $\sum \frac{(w_{ij}^{(\ell)})^2}{1 + (w_{ij}^{(\ell)})^2}$, what is $\frac{\partial \text{regularizer}}{\partial w_{ij}^{(\ell)}}$?

- ① $2w_{ij}^{(\ell)} / \left(1 + (w_{ij}^{(\ell)})^2 \right)^1$
- ② $2w_{ij}^{(\ell)} / \left(1 + (w_{ij}^{(\ell)})^2 \right)^2$
- ③ $2w_{ij}^{(\ell)} / \left(1 + (w_{ij}^{(\ell)})^2 \right)^3$
- ④ $2w_{ij}^{(\ell)} / \left(1 + (w_{ij}^{(\ell)})^2 \right)^4$

Reference Answer: ②

Too much calculus in this class, huh? :-)