# 1.2.16 Lesson Review

**Date:** 2/24/2026, 11:31:59 PM

**Time Spent:** 13:04

**Score: 100%**

Passing Score: 80%

## Question 1

Why are gradient-boosted decision trees well suited for tasks like ranking the likelihood that a newly registered domain is phishing in a supervised cybersecurity model?

- They rely exclusively on image-like representations of traffic so they can detect byte-level entropy patterns.

- They require no labeled data and instead learn normal behavior entirely from unlabeled traffic logs.

- They are restricted to modeling only linear relationships between a few numerical security features.

- They can integrate many fine-grained features and learn complex, non-linear relationships among them.  ✓ Correct

**Explanation**

Integrating many fine-grained features and learning complex, non-linear relationships among them is the correct answer. Gradient-boosted decision trees handle many heterogeneous features (domain age, registration, ASN, reputation, login failures) and capture non linear interactions, making them effective for supervised phishing-likelihood scoring.

Requiring no labeled data and instead learning normal behavior entirely from unlabeled traffic logs is incorrect. Gradient-boosted trees are supervised and need labeled examples (phishing vs. benign); unsupervised learning uses methods like clustering or Isolation Forests.

Relying exclusively on image-like representations of traffic so they can detect byte-level entropy patterns is incorrect. That describes convolutional neural networks on raw bytes, not boosted trees, which use engineered features.

Being restricted to modeling only linear relationships between a few numerical security features is incorrect. That fits logistic regression, not gradient-boosted trees, which model complex, non-linear interactions.

**Related Content**

📄 1.1.5 Machine Learning and Statistical Learning

📄 1.2.2 Supervised Learning

📄 1.2.3 Unsupervised Learning

resources\questions\q_supervised_learning_01.question.xml

A security architect is advising a healthcare organization that wants to use patient data from mobile devices to improve its diagnostic AI model, but strict privacy laws prohibit uploading raw patient data to the cloud. Which approach should the architect recommend to train the model while preserving data locality and compliance?

Use transfer learning by pretraining the model on public cloud datasets and ignoring local device data

Use centralized learning so all device data is encrypted and uploaded to a single training server

Use differential privacy to randomly share subsets of raw data from each device with the central server

● Use federated learning so the model is sent to devices, trained locally, and only model updates are returned ✓ Correct

**Explanation**

Using federated learning so the model is sent to devices, trained locally, and only model updates are returned is correct answer. Federated learning sends the global model to client devices, trains on local data that never leaves the device, and returns only parameters/updates, which supports privacy and regulatory compliance.

Using centralized learning so all device data is encrypted and uploaded to a single training server is incorrect. Centralized learning still requires aggregating raw data on a single server, violating the requirement that data stay on devices.

Using differential privacy to randomly share subsets of raw data from each device with the central server is incorrect. With differential privacy, raw or partially perturbed data is still shared with the server, conflicting with the stated constraints.

Using transfer learning by pretraining the model on public cloud datasets and ignoring local device data is incorrect. Transfer learning on public datasets ignores the valuable on-device data the organization explicitly wants to leverage.

**Related Content**

📄 1.2.5 Federated Learning

resources\questions\q_federated_learning_01.question.xml

**Question 3**                                        ⊘ **Correct**

When adapting a reinforcement learning–based security agent to perform better on a specific environment, what is the term for the phase in which the existing model is trained for additional epochs on new, targeted data?

    Hyperparameter search

    Fixed response playbooks

● Fine-tuning    ✓   Correct

    Model validation

**Explanation**

Fine-tuning is the correct answer. It continues training an already trained model on a narrower dataset for a few more epochs so it adapts to a specific task or environment. In security, this lets a reinforcement learning agent refine its policy for a given network, threat profile, or constraint while reusing prior learning.

Model validation is incorrect. It evaluates performance on unseen data to check generalization and overfitting and does not update model parameters.

Fixed response playbooks is incorrect. These are static incident-handling procedures, not a training phase for updating a model.

Hyperparameter search is incorrect. It optimizes settings like learning rate or batch size, not the actual fine-tuning of a pretrained model on targeted data.

**Related Content**

📄   1.2.4 Reinforcement Learning

resources\questions\q_reinforcement_learning_04.question.xml

Which outcomes describe key benefits of having a language model return results in a structured format such as JSON during security operations? (Select two.)

☑ Alerts can be enriched, prioritized, and routed automatically because fields like source IP and risk score    ✓    Correct
are predictable.

☐ The need for log normalization and field mapping in SIEM platforms is eliminated for all other data sources.

☑ Security tools can ingest and process the model's output with fewer custom parsers and less brittle integration    ✓    Correct
code.

☐ The model can directly decrypt captured network traffic without needing access to encryption keys or security appliances.

☐ Human analysts no longer need to review high-risk alerts because the structured format guarantees complete accuracy.

**Explanation**

Security tools can ingest and process the model's output with fewer custom parsers and less brittle integration code is a correct answer. Consistent JSON fields let SIEM, SOAR, and ticketing tools use simple ingestion rules instead of ad hoc text parsing.

Alerts can be enriched, prioritized, and routed automatically because fields like source IP and risk score are predictable is a correct answer. A stable schema lets downstream pipelines attach threat-intelligence data, sort by priority, and trigger playbooks based on specific keys, which speeds triage and response.

The model can directly decrypt captured network traffic without needing access to encryption keys or security appliances is incorrect. Structured output does not change cryptographic limits.

The need for log normalization and field mapping in SIEM platforms is eliminated for all other data sources is incorrect. Structured model output simplifies one stream, but other log sources still require normalization and mapping to a common schema.

Human analysts no longer need to review high risk alerts because the structured format guarantees complete accuracy is incorrect. Structured data improves consistency and automation, but does not ensure the model's judgments are always correct, so human oversight remains necessary.

**Related Content**

resources\questions\q_introduction_to_prompt_engineering_03.quesiton.question.xml

A security team is training a reinforcement learning–based firewall agent in a lab that replays mixed normal and attack traffic. The goal is for the firewall to automatically tune its rules so it blocks malicious sessions quickly while minimizing user impact.

Early in testing, the agent frequently rate-limits suspicious flows for long periods, even after those flows are confirmed malicious, causing noticeable latency for legitimate users on the same links.

What change BEST applies reinforcement learning to improve this adaptive firewall behavior?

Keep the reward table as-is but simplify the training data so it includes mostly one or two common attack patterns, reducing traffic variability so the agent converges faster, even if it becomes less flexible in handling new threats.

Disable the agent's ability to enforce rules and use it only for recommendations, while human analysts manually update firewall policies based on its suggestions and periodic review of incident reports.

Increase the reward for quickly blocking traffic later confirmed malicious and raise the penalty for prolonged throttling that causes latency, then retrain the agent on the same replayed traffic.   ✓ Correct

Train the agent in an environment that contains only benign sessions so it learns to avoid blocking and throttling, prioritizing user experience and stability over security decisions during its initial learning phase.

**Explanation**

Increasing the reward for quickly blocking traffic, later confirmed malicious, and raising the penalty for prolonged throttling that causes latency, then retraining the agent on the same replayed traffic is the correct answer. updates the reward to favor fast blocking and less throttling and re-runs the RL loop so the policy shifts toward more decisive blocking.

Disabling the agent's ability to enforce rules and using it only for recommendations is incorrect. This stops reward-driven policy updates.

Keeping the reward table as-is but simplifying training data to a few attack patterns is incorrect. It keeps the same incentives and reduces robustness.

Training the agent only on benign sessions is incorrect. It provides almost no feedback on malicious flows, so attack handling cannot improve.

**Related Content**

📄 1.2.4 Reinforcement Learning

resources\questions\q_reinforcement_learning_07.question.xml

In an unsupervised anomaly-detection system using autoencoders, what does a high reconstruction error for a particular log event usually indicate?

The event has been routed through multiple hidden layers to improve its encryption strength.

The event has been perfectly compressed and restored without any information loss.

The event was labeled as malicious during supervised training and is now ignored.

● The event differs significantly from the normal patterns the autoencoder learned.　✓ Correct

**Explanation**

The event differs significantly from the normal patterns the autoencoder learned is the correct answer. An autoencoder is trained to reconstruct "normal" records; if it cannot, the reconstruction error is high, indicating the event lies outside the learned baseline and may be anomalous.

The event has been perfectly compressed and restored without any information loss is incorrect. High reconstruction error means poor reconstruction, not perfect compression and restoration.

The event was labeled as malicious during supervised training and is now ignored is incorrect. Unsupervised autoencoders learn from unlabeled data and do not use malicious/benign labels.

The event has been routed through multiple hidden layers to improve its encryption strength is incorrect. Hidden layers are for representation learning, not encryption; anomaly signals come from reconstruction error.

**Related Content**

📄 1.1.5 Machine Learning and Statistical Learning

📄 1.2.3 Unsupervised Learning

resources\questions\q_unsupervised_learning_02.question.xml

What does it mean when a security model maintains high accuracy throughout adversarial stress tests?

It has passed a demanding validation stage, indicating it can be promoted directly from lab testing to full autonomous operation without the need for a shadow-mode comparison against human analysts.

It has demonstrated robust performance during evaluation, allowing teams to treat it as effectively immune to future data-poisoning attempts and to relax ongoing monitoring requirements.

Its decision logic can be trusted as resilient enough to support automated detection and response with rigor comparable to formal cryptographic key-ceremony procedures.

✓ Correct

It has shown strong generalization in testing, which justifies placing less emphasis on strict dataset separation and relying more on production feedback for future assessments.

**Explanation**

Its decision logic can be trusted as resilient enough to support automated detection and response with rigor comparable to formal cryptographic key ceremony procedures is the correct answer. High accuracy under adversarial stress tests indicates decision logic that is hard to manipulate and reliable enough for automated detection and response, with rigor comparable to cryptographic key ceremonies.

It has demonstrated robust performance during evaluation is incorrect. Strong stress-test performance is positive, but the model can still be affected by future data poisoning, concept drift, and evolving attacker tactics, so monitoring and validation cannot be relaxed.

It has shown strong generalization in testing is incorrect. Strict separation of training, validation, and test sets is a core protection against overfitting and silent data-poisoning. Strong stress-test results do not justify weakening this separation or relying mainly on production feedback, which would undermine unbiased evaluation.

It has passed a demanding validation stage is incorrect. Passing stress tests is not enough to skip cautious deployment steps like running in shadow mode before full autonomy.

**Related Content**

resources\questions\q_ai_model_training_02.question.xml

## Question 8

Why is exhaustive audit logging important when integrating a language model into security operations?

It encrypts all user credentials and tokens stored in the authentication system so that unauthorized access to the model management dashboard is prevented.

It automatically deletes all previous interactions after a short period, ensuring that no historical data exists for analysts or auditors to examine.

It provides a detailed record of every prompt and response, enabling investigators to reconstruct suspicious activity, trace misuse, and demonstrate compliance during reviews.    ✓ Correct

It continuously optimizes model hyperparameters based on logged performance metrics, improving classification accuracy without requiring a separate tuning workflow.

**Explanation**

Providing a detailed record of every prompt and response, enabling investigators to reconstruct suspicious activity, trace misuse, and demonstrate compliance during reviews, is correct. Exhaustive audit logging records timestamps, user IDs, IPs, prompts, and responses so teams can replay incidents, spot abuse patterns, and provide evidence for audits.

Encrypting all user credentials and tokens stored in the authentication system is incorrect. That is an identity, encryption, and secure-storage function, not logging.

Automatically deleting all previous interactions after a short period is incorrect. It undermines the value of persistent records for investigation and compliance.

Continuously optimizing model hyperparameters based on logged metrics is incorrect. Logging provides visibility and evidence; it does not itself tune or reconfigure the model.

**Related Content**

resources\questions\q_securing_the_model_03.question.xml

In designing user prompts for security workflows, what is the primary reason to explicitly specify the desired output format (for example, particular JSON keys)?

It guarantees that downstream parsers, orchestration tools, and reporting dashboards can consume the model's output consistently.

● It allows downstream tools and scripts to reliably parse and act on the model's response.                                  ✓  Correct

It guarantees that the model will always reach the correct security conclusion.

It prevents the model from generating any natural-language explanation.

**Explanation**

Allowing downstream tools and scripts to reliably parse and act on the model's response is the correct answer. Defining an explicit output format makes responses predictable and machine-readable (for example, via JSON keys), so automation can reliably parse and use them in security workflows.

Guaranteeing that downstream parsers, orchestration tools, and reporting dashboards can consume the model's output consistently is incorrect. Structured formats improve consistency but do not guarantee flawless integration or remove the need for validation.

Guaranteeing that the model will always reach the correct security conclusion is incorrect. Format affects structure, not correctness.

Preventing the model from generating any natural-language explanation is incorrect. You can request both structured fields for tools and narrative explanations for human analysts.

**Related Content**

📄  1.2.12 User Prompts

resources\questions\q_user_prompts_01.question.xml

An incident responder wants to quickly assess whether a new, suspicious DNS query might indicate malicious activity.

Historical labeled examples are not immediately available, and the responder needs an initial risk assessment from an AI assistant without spending time on data preparation.

Which approach BEST applies zero-shot prompting in this situation?

- Provide the raw DNS query to the AI with a single clear instruction such as, "Analyze this DNS query and assess its likelihood of being malicious on a scale from 1–10, explaining your reasoning," without supplying any prior labeled examples.                    ✓   Correct

  Construct a template that automatically inserts the new DNS query alongside several known benign queries, then ask the AI to identify which query is most likely malicious, returning the result in structured JSON for ingestion by existing dashboards.

  Build a reusable template that injects multiple days of mixed DNS logs and their analyst-assigned labels, then ask the AI to learn from these examples and generate rules for automatically categorizing future queries by risk level.

  Paste the DNS query into a standardized prompt template that includes ten labeled examples of past DNS queries, half benign and half malicious, and ask the AI to classify the new query based on those examples and explain its reasoning.

**Explanation**

Provide the raw DNS query to the AI with a single clear instruction such as, "Analyze this DNS query and assess its likelihood of being malicious on a scale from 1–10, explaining your reasoning," without supplying any prior labeled examples is the correct answer. It is true zero-shot: only the raw query plus a task description, no examples.

Paste the DNS query into a standardized prompt template that includes ten labeled examples of past DNS queries is incorrect. This is multi-shot prompting, not zero-shot.

Construct a template that automatically inserts the new DNS query alongside several known benign queries is incorrect. Those benign examples make it one/few-shot, not zero-shot.

Build a reusable template that injects multiple days of mixed DNS logs and their analyst-assigned labels is incorrect. This is many-shot, example-driven prompting, not zero-shot.

**Related Content**

📄  1.2.13 Zero-Shot, One-Shot, Multi-Shot, and Templates
resources\questions\q_zero-shot_one-shot_multi-shot_and_templates_02.question.xml