

OWASP Makine Öğrenmesi Güvenliđi İlk 10 – Taslak Sürüm v0.3



İçindekiler

Giriş	5
! Önemli	5
Genel Bakış	5
Hedef Kitle	5
Kapsam	5
Dikkat	6
Yayın.....	6
Ana Yazarlar	6
Telif Hakkı ve Lisans.....	6
OWASP Hakkında.....	7
İlk 10 2023 Listesi	8
ML01:2023 Girdi (Input) Manipülasyon Saldırıları.....	9
Tanım.....	9
Nasıl Önlenir	9
Risk Faktörleri.....	9
Örnek Saldırı Senaryoları:.....	10
Senaryo #1: Görüntü Sınıflandırma sistemlerinin girdi (input) manipülasyonu.....	10
Senaryo #2: Saldırı Tespit Sistemlerinden kaçmak için ağ trafiğinin manipülasyonu.....	10
Referanslar	10
ML02:2023 Veri Zehirlenmesi Saldırısı	11
Tanım:.....	11
Nasıl Önlenir:	11
Risk Faktörleri.....	11
Örnek Saldırı Senaryoları:.....	12
Senaryo #1: Spam sınıflandırıcısı eğitimi	12
Senaryo #2: Ağ trafiği sınıflandırma sisteminin eğitilmesi.....	12
Referanslar	12
ML03:2023 Model Ters Çevirme Saldırısı	13
Tanım.....	13
Nasıl Önlenir:	13
Risk Faktörleri.....	13
Saldırı Senaryosu Örneği	14
Senaryo #1: Bir yüz tanıma modelinden kişisel bilgilerin çalınması.....	14
Senaryo #2: Çevrimiçi reklamcılıkta bir bot tespit modelini atlama	14
Referanslar	14



ML04:2023 Üyelik Çıkarım Saldırısı	15
Tanım	15
Nasıl Önlenir	15
Risk Faktörleri	15
Örnek Saldırı Senaryoları:.....	16
Senaryo #1: Bir makine öğrenimi modelinden finansal verileri çıkarma	16
Referanslar	16
ML05:2023 Model Çalma.....	17
Tanım:	17
Nasıl Önlenir:	17
Risk Faktörleri	17
Saldırı Senaryosu Örneği:	18
Senaryo #1: Bir rakipten makine öğrenimi modeli çalmak	18
Referanslar	18
ML06:2023 Yapay Zeka Tedarik Zinciri Saldırıları.....	19
Tanım:	19
Nasıl Önlenir:	19
Risk Faktörleri	19
Saldırı Senaryosu Örneği	20
Senaryo #1: Bir organizasyondaki makine öğrenimi projesine saldırı.....	20
Referanslar	20
ML07:2023 Transfer Öğrenme Saldırısı.....	21
Tanım	21
Nasıl Önlenir	21
Risk Faktörleri	21
Saldırı Senaryosu Örneği	22
Senaryo #1: Bir modelin kötü niyetli bir veri seti üzerinde eğitilmesi.....	22
Referanslar	22
ML08:2023 Model Eğriliği	23
Tanım	23
Nasıl Önlenir	23
Risk Faktörleri	24
Saldırı Senaryosu Örneği:	24
Senaryo #1: Model eğriliği ile finansal kazanç	24
Referanslar	24
ML09:2023 Çıktı Bütünlüğü Saldırısı	25
Tanım:	25



Nasıl Önlenir:	25
Risk Faktörleri	25
Saldırı Senaryosu Örneği	26
Senaryo #1: Hasta sağlık kayıtlarının değiştirilmesi.....	26
Referanslar	26
ML10:2023 Model Zehirlenme	27
Tanım	27
Nasıl Önlenir	27
Risk Faktörleri	27
Saldırı Senaryosu Örneği	28
Senaryo #1: Model zehirlenme yoluyla finansal kazanç elde etme	28
Referanslar	28
Teşekkürler	29
Katkıda Bulunanlar	29
Nasıl katkıda bulunulur	29
Terimler Sözlüğü (Tamamlanmamış)	30



Giriş

! Önemli

Bu çalışmanın güncel versiyonu taslak halindedir ve sık sık değiştirilmektedir. Projenin yayınlanma takvimi hakkında ve projeye nasıl katkıda bulunulacağı hakkında bilgi için lütfen proje wiki'sine bakın.

Genel Bakış

OWASP Makine Öğrenimi Güvenliği İlk 10 projesinin temel amacı, makine öğrenimi (ML) sistemlerinin (en önemli ilk 10 güvenlik sorununu genel olarak ele almaktır. Bu nedenle, bu projenin başlıca hedeflerinden biri, endüstri meslektaşları tarafından incelenen yüksek kaliteli bir sonuç üretmektir.

Hedef Kitle

Bu proje yayımlanırken başlıca hedeflenen kitle; geliştiriciler, makine öğrenimi mühendisleri (MLE) ve operasyon uygulayıcıları ile uygulama güvenliği uzmanlarıdır. Bu rollerin her biri makine öğrenimi sistemleri oluşturur, işletir ve güvence altına alırken, içerik herkese hitap edilecek şekilde tasarlanmıştır. İçerik, uygun olduğu durumlarda belirli teknoloji alanları için gereken anlayış düzeyini belirlemeyi amaçlayacaktır.

Kapsam

Bu proje, makine öğrenimi (ML) sistemlerinin (en önemli ilk 10 güvenlik sorununun genel bir bakışını sunacaktır. Makine öğrenimi sistemlerinin hızlı benimsenmesi nedeniyle, bu projeden daha dar veya daha geniş kapsama sahip olan ilgili projeler OWASP ve benzeri organizasyonlar içinde bulunabilir. Örnek olarak, düşmanca saldırılar (adversarial attacks) bir tehdit kategorisi olsa da, bu proje aynı zamanda makine öğreniminin operasyonel ve mühendislik süreçlerinde güvenlik gibi düşmanca olmayan senaryoları da ele alacaktır.



Dikkat

Yayın

Bu belge Őu anda v0.3 taslak sűrűműndedir.

- [Sűrűm Notları](#)
- [DeęiŐiklik Gűnlűęű](#)

Ana Yazarlar

- [Shain Singh](#)
- [Sagar Bhure](#)
- [Rob van der Veer](#)

Telif Hakkı ve Lisans



Telif Hakkı © 2003-2023 OWASP Vakfı'na aittir. Bu belge, Creative Commons Attribution Share-Alike 4.0 lisansı altında yayınlanmıŐtır. Herhangi bir yeniden kullanım veya daęıtım durumunda, bu alıŐmanın lisans Őartlarını dięer kiŐilere net bir Őekilde belirtmelisiniz."



OWASP Hakkında

Açık Web Uygulama Güvenliği Projesi (The Open Worldwide Application Security Project (OWASP)), organizasyonların uygulamalar ve API'lar geliştirmelerine, satın almalarına ve sürdürmelerine güvenebilecekleri bir olanak tanımaya adanmış açık bir topluluktur."

OWASP'ta ücretsiz ve açık kaynaklı olarak bulabileceğiniz şunlar vardır:

- Uygulama güvenliği araçları ve standartları.
- Uygulama güvenliği testi, güvenli kod geliştirme ve güvenli kod incelemesi konularındaki eksiksiz kitaplar.
- Sunumlar ve [videolar](#).
- Birçok yaygın konuya yönelik [özet bilgiler](#).
- Standart güvenlik kontrolleri ve kütüphaneler.
- [Dünya genelinde yerel şubeler](#).
- Son teknoloji araştırmalar.
- Dünya genelinde geniş kapsamlı [konferanslar](#).
- [E-posta listeleri \(arşivlenmiş\)](#).

Daha fazla bilgi için: <https://www.owasp.org>.

Tüm OWASP araçları, belgeleri, videoları, sunumları ve, şubeleri uygulama güvenliğini geliştirmek isteyen herkes için ücretsiz ve açıktır.

Uygulama güvenliğine yaklaşımı bir insan, süreç ve teknoloji sorunu olarak benimseriz, çünkü uygulama güvenliği konusundaki en etkili yaklaşımlar bu alanlarda gelişmeler gerektirir.

OWASP, yeni türde bir organizasyondur. Ticari baskılardan özgür olmamız, uygulama güvenliği hakkında tarafsız, pratik ve maliyet-etkin bilgi sağlayabilmemize olanak tanır.

OWASP, herhangi bir teknoloji şirketi ile ilişkili değildir, ancak ticari güvenlik teknolojilerinin bilinçli kullanımını destekleriz. OWASP, birçok türde materyalini işbirlikçi, şeffaf ve açık bir şekilde üretir.

OWASP Vakfı, projenin uzun vadeli başarısını sağlayan kar amacı gütmeyen bir kuruluştur. OWASP ile ilişkilendirilen neredeyse herkes, OWASP yönetim kurulu, şube liderleri, proje liderleri ve proje üyeleri dahil olmak üzere gönüllüdür.

Biz, araştırma hibeleri ve altyapı ile yenilikçi güvenlik araştırmalarımızı destekliyoruz.

Bize katılın!



İlk 10 2023 Listesi



ML01:2023 Girdi (Input) Manipülasyon Saldırıları

Tanım

Girdi Manipülasyon Saldırıları, bir saldırganın modeli yanıltmak için girdi verilerini kasıtlı olarak değiştirdiği bir saldırı türü olan Düşmanca (adversarial) Saldırıları içeren bir kapsayıcı terimdir.

Bu tablonun yalnızca aşağıdaki senaryoya dayalı bir örnek olduğunu belirtmek önemlidir. Gerçek risk değerlendirmesi, her bir makine öğrenimi sisteminin kendine özgü koşullarına bağlı olacaktır.

Nasıl Önlenir

Düşmanca Eğitim: Girdi manipülasyonu saldırısına yönelik savunmaya ilişkin yaklaşımlardan biri modeli düşmanca örnekler üzerinde eğitmektir. Bu, modelin saldırılara karşı daha dayanıklı olmasına ve yanlış yönlendirilmeye karşı duyarlılığını azaltmasına yardımcı olabilir.

Sağlam modeller: Diğer bir yaklaşım, düşmanca eğitim veya savunma mekanizmalarını içeren modeller gibi düşmanca saldırılara karşı sağlam olacak şekilde tasarlanmış modelleri kullanmaktır.

Girdi doğrulama: Girdi doğrulama, girdi manipülasyonu saldırılarını tespit etmek ve önlemek için kullanılacak bir diğer önemli savunma mekanizmasıdır. Bu, girdi verilerinin beklenmedik değerler veya kalıplar gibi anormalliklere karşı kontrol edilmesini ve kötü amaçlı olması muhtemel girdilerin reddedilmesini içerir.

Risk Faktörleri

Tehdit Aracıları/ Saldırı Vektörleri	Güvenlik Zafiyetleri	Etki
İstismar Edilebilirlik: 5 (kolay) ML Uygulama Özel: 4 ML İşlemlere Özel: 3	Tespit Edilebilirlik: 3 (Orta) Manipüle edilen görüntü çıplak gözle fark edilemeyebilir, bu da saldırının tespit edilmesini zorlaştırır.	Teknik: 5 (Zor) Saldırı, derin öğrenme ve görüntü işleme teknikleri hakkında teknik bilgi gerektirir.
Tehdit Aracıları: Derin öğrenme ve görüntü işleme teknikleri hakkında bilgi sahibi olan saldırgan. Saldırı Vektörü: Meşru bir görüntüye benzeyen, kasıtlı olarak hazırlanmış manipüle edilmiş görüntü.	Derin öğrenme modelinin görüntüleri doğru bir şekilde sınıflandırma becerisindeki güvenlik açığı. Derin öğrenme modelinin görüntüleri doğru bir şekilde sınıflandırma becerisindeki güvenlik açığı.	Görüntünün yanlış sınıflandırılması, güvenliğin atlanmasına veya sistemin zarar görmesine neden olur. Görüntünün yanlış sınıflandırılması, güvenliğin atlanmasına veya sistemin zarar görmesine neden olur.

Bu tablonun yalnızca [aşağıdaki senaryoya](#) dayalı bir örnek olduğunu unutmamak önemlidir. Gerçek risk değerlendirmesi, her bir makine öğrenimi sisteminin kendine özgü koşullarına bağlı olacaktır.



Örnek Saldırı Senaryoları

Senaryo #1: Görüntü Sınıflandırma sistemlerinin girdi (input) manipülasyonu

Derin öğrenme modeli, görüntüleri köpekler ve kediler gibi farklı kategorilere ayırmak için eğitilmiştir. Bir saldırgan, meşru bir kedi görüntüsüne çok benzeyen ancak modelin bunu bir köpek olarak yanlış sınıflandırmasına neden olabilecek küçük, dikkatlice hazırlanmış pertürbasyonlarla orijinal görüntüyü manipüle eder. Model gerçek dünya ortamına yerleştirildiğinde, saldırgan manipüle edilmiş görüntüyü güvenlik önlemlerini atlamak veya sisteme zarar vermek için kullanabilir.

Senaryo #2: Saldırı Tespit Sistemlerinden kaçmak için ağ trafiğinin manipülasyonu

Bir derin öğrenme modeli, bir ağdaki izinsiz girişleri tespit etmek için eğitilir. Bir saldırgan, modelin izinsiz giriş tespit sisteminden kaçacak şekilde paketleri dikkatlice hazırlayarak ağ trafiğini manipüle eder. Saldırgan, kaynak IP adresi, hedef IP adresi veya yük gibi ağ trafiğinin özelliklerini, saldırı tespit sistemi tarafından algılanmayacak şekilde değiştirebilir. Örneğin, saldırgan kaynak IP adresini bir proxy sunucusunun arkasına gizleyebilir veya ağ trafiğinin yükünü şifreleyebilir. Bu tür bir saldırı, veri hırsızlığına, sistemin ele geçirilmesine veya diğer hasar biçimlerine yol açabileceğinden ciddi sonuçlar doğurabilir.

Referanslar



ML02:2023 Veri Zehirlenmesi Saldırısı

Tanım

Veri zehirlenmesi saldırıları, bir saldırgan modelin istenmeyen bir şekilde davranmasına neden olmak için eğitim verilerini değiştirdiğinde oluşur.

Nasıl Önlenir

Veri doğrulama ve onaylama: Modeli eğitmede kullanılmadan önce eğitim verilerinin kapsamlı bir şekilde doğrulandığından ve onaylandığından emin olun. Bunu, veri doğrulama kontrolleri uygulayarak ve veri etiketlemenin doğruluğunu onaylamak için birden fazla veri etiketleyiciyi kullanarak yapabilirsiniz.

Güvenli veri depolama: Eğitim verilerini şifreleme, güvenli veri aktarım protokolleri ve güvenlik duvarları kullanma gibi yöntemlerle güvenli bir şekilde depolayın.

Veri ayırma: Eğitim verilerinin tehlikeye girme riskini azaltmak için eğitim verilerini üretim verilerinden ayırın.

Erişim denetimi: Eğitim verilerine kimlerin erişebileceğini ve ne zaman erişebileceklerini sınırlamak için erişim kontrolleri uygulayın.

İzleme ve denetim: Herhangi bir anormallik için eğitim verilerini düzenli olarak izleyin ve herhangi bir veri tahrifatını tespit etmek için denetimler yapın.

Model doğrulama: Modeli, eğitim sırasında kullanılmamış ayrı bir doğrulama seti kullanarak doğrulayın. Bu, eğitim verilerini etkilemiş olabilecek veri zehirlenme saldırılarının tespit edilmesine yardımcı olabilir.

Model toplulukları: Eğitim verilerinin farklı alt kümelerini kullanarak birden fazla model eğitin ve tahminler yapmak için bu modellerden oluşan bir topluluk kullanın. Bu, saldırganın hedeflerine ulaşmak için birden fazla model tehlikeye atmasını gerektireceğinden veri zehirlenme saldırılarının etkisini azaltabilir.

Anomali tespiti: Eğitim verilerindeki, örneğin veri dağılımında veya veri etiketlemesinde ani değişiklikler gibi anormal davranışları tespit etmek için anomali tespit tekniklerini kullanın. Bu teknikler veri zehirlenme saldırılarını erkenden tespit etmek için kullanılabilir.

Risk Faktörleri

Tehdit Aracıları/ Saldırı Vektörleri	Güvenlik Zafiyetleri	Etki
Sömürülebilirlik: 3 (Orta Seviye) ML Uygulama Özel: 4 ML İşlemlere Özel: 3	Tespit Edilebilirlik: 2 (Zor)	Teknik: 4
Tehdit Aracısı: Model için kullanılan eğitim verilerine erişimi olan saldırgan. Saldırı Vektörü: Saldırgan, eğitim veri kümesine kötü amaçlı veriler ekler.	Veri doğrulama eksikliği ve eğitim verilerinin yetersiz izlenmesi.	Model, zehirlenmiş verilere dayanarak yanlış tahminlerde bulunacak, yanlış kararlara ve ciddi sonuçlara yol açacaktır.



Bu tablonun yalnızca [aşağıdaki senaryoya](#) dayalı bir örnek olduğunu unutmamak önemlidir. Gerçek risk değerlendirmesi, her bir makine öğrenimi sisteminin kendine özgü koşullarına bağlı olacaktır.

Örnek Saldırı Senaryoları

Senaryo #1: Spam sınıflandırıcısı eğitimi

Bir saldırgan, e-postaları spam ya da spam olmayan olarak sınıflandıran bir derin öğrenme modelinin eğitim verilerini zehirliyor. Saldırgan bu saldırıyı, kötü niyetli olarak etiketlenmiş spam e-postaları eğitim veri setine enjekte ederek gerçekleştirir. Bu, örneğin ağa girerek veya veri depolama yazılımındaki bir güvenlik açığından yararlanarak veri depolama sistemini ele geçirebilir. Saldırgan ayrıca e-posta etiketlemede değişiklik yaparak ya da veri etiketleme işini yapan kişilere rüşvet vererek yanlış etiketleme yapmalarını sağlamak gibi yöntemlerle veri etiketleme sürecini manipüle edebilir.

Senaryo #2: Ağ trafiği sınıflandırma sisteminin eğitilmesi

Bir saldırgan, ağ trafiğini e-posta, web gezintisi ve video akışı gibi farklı kategorilere sınıflandırmak için kullanılan bir derin öğrenme modelinin eğitim verilerini zehirler. Farklı bir trafik türü olarak yanlış etiketlenmiş çok sayıda ağ trafiği örneği sunarak modelin bu trafiği yanlış kategori olarak sınıflandırmak üzere eğitilmesine sebep olur. Sonuç olarak, model dağıtıldığında yanlış trafik sınıflandırmaları yapmak üzere eğitilebilir ve bu da potansiyel olarak ağ kaynaklarının yanlış tahsisine veya ağ performansının düşmesine yol açabilir.

Referanslar



ML03:2023 Model Ters Çevirme Saldırısı

Tanım

Model ters çevirme saldırıları, bir saldırganın modelden bilgi çıkarmak için tersine mühendislik yaptığı durumlarda ortaya çıkar.

Nasıl Önlenir

Erişim kontrolü: Modele veya tahminlerine erişimin sınırlandırılması, saldırganların modeli tersine çevirmek için gereken bilgileri elde etmesini önleyebilir. Bu, modele veya tahminlerine erişirken kimlik doğrulama, şifreleme veya diğer güvenlik biçimlerini gerekli kılarak yapılabilir.

Girdi doğrulama: Modelin girdilerinin doğrulanması, saldırganların modeli tersine çevirmek için kullanılacak kötü niyetli veriler sağlamasını önleyebilir. Bu, model tarafından işlenmeden önce girdilerin formatını, aralığını ve tutarlılığını kontrol ederek yapılabilir.

Model şeffaflığı: Modeli ve tahminlerini şeffaf hale getirmek, model tersine çevirme saldırılarını tespit etmeye ve önlemeye yardımcı olabilir. Bu, tüm girdileri ve çıktıları loglayarak, modelin tahminleri için açıklamalar sağlayarak veya kullanıcıların modelin dahili temsillerini denetlemesine izin vererek yapılabilir.

Düzenli izleme: Modelin anomali tahminlerinin izlenmesi model ters çevirme saldırılarının tespit edilmesine ve önlenmesine yardımcı olabilir. Bu, girdilerin ve çıktıların dağılımını izleyerek, modelin tahminlerini temel gerçek verilerle karşılaştırarak veya modelin zaman içindeki performansını izleyerek yapılabilir.

Modelin yeniden eğitilmesi: Modelin düzenli olarak yeniden eğitilmesi, model ters çevirme saldırıları tarafından sızdırılan bilgilerin güncelliğini yitirmesine engel olabilir. Bu, yeni verileri dahil ederek ve modelin tahminlerindeki yanlışlıkları düzelterek yapılabilir.

Risk Faktörleri

Tehdit Aracıları/ Saldırı Vektörleri	Güvenlik Zafiyetleri	Etki
Sömürülebilirlik: 4 (Orta) ML Uygulama Özel: 5 ML İşlemlere Özel: 3	Tespit Edilebilirlik: 2 (Zor)	Teknik: 4 (Orta düzeyde teknik beceri gereklidir)
Tehdit Aracıları: Modele ve girdi verilerine erişimi olan saldırganlar. Saldırı Vektörleri: Modele bir görüntü gönderme ve modelin yanıtını analiz etme	Modelin çıktısı, girdi verileri hakkında hassas bilgiler çıkarmak için kullanılabilir.	Girdi verileriyle ilgili gizli bilgiler tehlikeye girebilir.

Bu tablonun yalnızca [aşağıdaki senaryoya](#) dayalı bir örnek olduğunu unutmamak önemlidir. Gerçek risk değerlendirmesi, her bir makine öğrenimi sisteminin kendine özgü koşullarına bağlı olacaktır.



Saldırı Senaryosu Örneği

Senaryo #1: Bir yüz tanıma modelinden kişisel bilgilerin çalınması

Saldırgan, yüz tanıma işlemini gerçekleştirmek için bir derin öğrenme modelini eğitir. Daha sonra bu modeli, bir şirket veya kuruluş tarafından kullanılan farklı bir yüz tanıma modeline model ters çevirme saldırısı gerçekleştirmek için kullanır. Saldırgan, bireylerin görüntülerini modele girer ve modelin tahminlerinden bireylerin adları, adresleri veya sosyal güvenlik numaraları gibi kişisel bilgilerini kurtarır.

Saldırgan bu saldırıyı, modeli yüz tanıma gerçekleştirecek şekilde eğiterek ve ardından bu modeli başka bir yüz tanıma modelinin tahminlerini tersine çevirmek için kullanarak gerçekleştirmiştir. Bu, modelin hayata geçirilmesindeki bir güvenlik açığından yararlanılarak ya da modele bir API aracılığıyla erişilerek yapılabilir. Saldırgan daha sonra modelin tahminlerinden bireylerin kişisel bilgilerini kurtarabilir.

Senaryo #2: Çevrimiçi reklamcılıkta bir bot tespit modelini atlama

Bir reklamveren, reklamlara tıklama ve web sitelerini ziyaret etme gibi eylemleri gerçekleştirmek için botları kullanarak reklam kampanyalarını otomatikleştirmek ister. Ancak, çevrimiçi reklam platformları botların bu eylemleri gerçekleştirmesini önlemek için bot algılama modelleri kullanır. Reklamveren bu modelleri atlatmak için bot tespiti amacıyla bir derin öğrenme modeli eğitir ve bunu çevrimiçi reklam platformu tarafından kullanılan bot tespit modelinin tahminlerini tersine çevirmek için kullanır. Reklamveren botlarını modele girer ve botların insan kullanıcılar gibi görünmesini sağlayarak bot tespitini atlatabilir ve otomatik reklam kampanyalarını başarıyla yürütebilir.

Reklamveren bu saldırıyı, kendi bot tespit modelini eğiterek ve daha sonra bunu çevrimiçi reklam platformu tarafından kullanılan bot tespit modelinin tahminlerini tersine çevirmek için kullanarak gerçekleştirmiştir. Diğer modele, modelin hayata geçirilmesindeki bir açık üzerinden ya da bir API kullanarak ulaşmışlardır. Saldırının nihai sonucu, reklamverenin botlarını insan kullanıcı gibi görünmesini sağlayarak reklam kampanyalarını başarılı bir şekilde otomatikleştirmesidir.

Referanslar



ML04:2023 Üyelik Çıkarım Saldırısı

Tanım

Üyelik çıkarımı saldırıları, saldırganın modelin eğitim verilerini manipüle ederek hassas bilgileri açığa çıkaracak şekilde davranmasını sağlamasıyla gerçekleşir.

Nasıl Önlenir

Rastgele veya karıştırılmış veriler üzerinde model eğitimi: Makine öğrenimi modellerinin rastgele veya karıştırılmış veriler üzerinde eğitilmesi, bir saldırganın belirli bir örneğin eğitim veri kümesine dahil edilip edilmediğini belirlemesini daha zor hale getirebilir.

Model Gizleme: Rastgele gürültü ekleyerek veya farklı gizlilik tekniklerini kullanarak modelin tahminlerini gizlemek, bir saldırganın modelin eğitim verilerini belirlemesini zorlaştırarak üyelik çıkarımı saldırılarının önlenmesine yardımcı olabilir.

Düzenleştirme: L1 veya L2 düzenleştirme gibi düzenleştirme teknikleri, modelin eğitim verilerine aşırı uyumunu önlemeye yardımcı olabilir, bu da modelin belirli bir örneğin eğitim veri kümesine dahil edilip edilmediğini doğru bir şekilde belirleme yeteneğini azaltabilir.

Eğitim verilerini azaltma: Eğitim veri kümesinin boyutunu azaltmak veya gereksiz ya da yüksek korelasyonlu özellikleri kaldırmak, bir saldırganın üyelik çıkarımı saldırısından elde edebileceği bilgileri azaltmaya yardımcı olabilir.

Test etme ve izleme: Modelin davranışını anormalliklere karşı düzenli olarak test etmek ve izlemek, bir saldırganın hassas bilgilere erişmeye çalıştığını saptayıp üyelik çıkarımı saldırılarını tespit etmeye ve önlemeye yardımcı olabilir.

Risk Faktörleri

Tehdit Aracıları/ Saldırı Vektörleri	Güvenlik Zafiyetleri	Etki
İstismar Edilebilirlik: 4 (Orta) ML Uygulama Özel: 5 ML İşlemlere Özel: 3	Tespit Edilebilirlik: 3 (Orta)	Teknik: 4 (Orta)
Tehdit Aracıları: Verilere ve modele erişimi olan bilgisayar korsanları veya kötü niyetli aktörler. Verilere müdahale etmek için kötü niyetli olan veya rüşvet alan içeriden kişiler. Saldırı Vektörleri: Verilere yetkisiz erişime izin veren güvenli olmayan veri iletim kanalları.	Uygun veri erişim kontrollerinin eksikliği. Uygun veri doğrulama ve temizleme tekniklerinin eksikliği. Uygun veri şifreleme eksikliği. Uygun veri yedekleme ve kurtarma tekniklerinin eksikliği.	Güvenilmez veya yanlış model tahminleri. Hassas verilerin gizliliğinin ve mahremiyetinin kaybı. Yasal ve düzenleyici uyumluluk ihlalleri. İtibar kaybı.

Bu tablonun yalnızca [aşağıdaki senaryoya](#) dayalı bir örnek olduğunu unutmamak önemlidir. Gerçek risk değerlendirmesi, her bir makine öğrenimi sisteminin kendine özgü koşullarına bağlı olacaktır.



Örnek Saldırı Senaryoları:

Senaryo #1: Bir makine öğrenimi modelinden finansal verileri çıkarma

Kötü niyetli bir saldırgan, bireylerin hassas finansal bilgilerine erişmek ister. Bunu, bir makine öğrenimi modelini mali kayıtlardan oluşan bir veri kümesinde eğiterek ve belirli bir bireyin kaydının eğitim verilerine dahil edilip edilmediğini sorgulamak için kullanarak yaparlar. Saldırgan daha sonra bu bilgileri, bireylerin mali geçmişini ve hassas bilgilerini çıkarmak için kullanabilir.

Saldırgan daha sonra bu bilgiyi bireylerin finansal geçmişini ve hassas bilgilerini çıkarmak için kullanabilir. Saldırgan bu saldırıyı, bir finans kuruluşundan elde edilen mali kayıtlardan oluşan bir veri kümesi üzerinde makine öğrenimi modelini eğiterek gerçekleştirmiştir. Daha sonra bu modeli, belirli bir kişinin kaydının eğitim verilerine dahil edilip edilmediğini sorgulamak için kullanıp hassas finansal bilgileri ortaya çıkarmışlardır.

Referanslar



ML05:2023 Model Çalma

Tanım:

Model çalma saldırıları, bir saldırgan modelin parametrelerine erişim kazandığında meydana gelir.

Nasıl Önlenir

Şifreleme: Modelin kodunun, eğitim verilerinin ve diğer hassas bilgilerin şifrenmesi, saldırganların modele erişmesini ve modeli çalmasını önleyebilir.

Erişim Kontrolü: İki faktörlü kimlik doğrulama gibi sıkı erişim kontrol önlemlerinin uygulanması, yetkisiz kişilerin modele erişmesini ve çalmasını önleyebilir.

Düzenli yedeklemeler: Modelin kodunu, eğitim verilerini ve diğer hassas bilgilerini düzenli olarak yedeklemek, bir hırsızlık durumunda kurtarılmasını sağlayabilir.

Model Gizleme: Modelin kodunu gizlemek ve tersine mühendislik yapmayı zorlaştırmak, saldırganların modeli çalmasını engelleyebilir.

Filigranlama: Modelin koduna ve eğitim verilerine filigran eklemek, bir hırsızlığın kaynağını izlemeyi ve saldırganı sorumlu tutmayı mümkün kılabilir.

Yasal koruma: Model için patentler veya ticari sırlar gibi yasal koruma sağlamak, bir saldırganın modeli çalmasını daha zor hale getirebilir ve bir hırsızlık durumunda yasal işlem için bir temel sağlayabilir.

İzleme ve denetim: Modelin kullanımının düzenli olarak izlenmesi ve denetlenmesi, bir saldırganın modele erişmeye veya modeli çalmaya çalıştığını tespit ederek hırsızlığın saptanmasına ve önlenmesine yardımcı olabilir.

Risk Faktörleri

Tehdit Aracıları/ Saldırı Vektörleri	Güvenlik Zafiyetleri	Etki
İstismar Edilebilirlik: 4 (Orta) ML Uygulamasına Özel: 4 ML İşlemlere Özel: 3	Tespit Edilebilirlik: 3 (Orta)	Teknik: 4 (Orta)
Aracı/Saldırı Vektörü: Bu, saldırıyı gerçekleştiren varlığı ifade eder, bu durumda, makine öğrenimi modelini çalmak isteyen bir saldırgandır.	Güvenli olmayan model kurulumu: Modelin güvenli olmayan bir şekilde kurulması, saldırganın modele erişmesini ve modeli çalmasını kolaylaştırır.	Bir model hırsızlığının etkisi hem modeli eğitmek için kullanılan verilerin gizliliği hem de modeli geliştiren kuruluşun itibarı üzerinde olabilir.

Bu tablonun yalnızca [aşağıdaki senaryoya](#) dayalı bir örnek olduğunu unutmamak önemlidir. Gerçek risk değerlendirmesi, her bir makine öğrenimi sisteminin kendine özgü koşullarına bağlı olacaktır.



Saldırı Senaryosu Örneđi:

Senaryo #1: Bir rakipten makine öğrenimi modeli çalmak

Kötü niyetli bir saldırgan, değerli bir makine öğrenimi modeli geliřtirmiş olan bir řirketin rakibi için çalışmaktadır. Saldırgan bu modeli çalmak ister, böylece řirket rekabet avantajı elde edebilir ve kendi amaçları için kullanmaya başlayabilir.

Saldırgan bu saldırıyı, binary kodu parçalarına ayırarak ya da modelin eğitim verilerine ve algoritmasına erişerek řirketin makine öğrenimi modeline tersine mühendislik uygulayarak gerçekleştirir. Saldırgan model üzerinde tersine mühendislik yaptıktan sonra bu bilgileri kullanarak modeli yeniden oluşturabilir ve kendi amaçları doğrultusunda kullanmaya başlayabilir. Bu, orijinal řirket için önemli mali kayıpların yanı sıra itibarının da hasar görmesine neden olabilir.

Referanslar



ML06:2023 Yapay Zeka Tedarik Zinciri Saldırıları

Tanım

"AI Tedarik Zinciri Saldırıları," bir saldırganın bir sistem tarafından kullanılan bir makine öğrenimi kütüphanesini veya modelini değiştirdiği veya yerine koyduğu zaman meydana gelir. Bu aynı zamanda makine öğrenimi modelleri ile ilişkilendirilen verileri içerebilir.

Nasıl Önlenir

Paket İmzalarını Doğrulayın: Herhangi bir paketi kurmadan önce, paketlerin dijital imzalarını doğrulayarak bunların değiştirilmediğinden emin olun.

Güvenli Paket Depolarını (Repositories) Kullanın: Anaconda gibi güvenli paket depolarını kullanın, bunlar sıkı güvenlik önlemleri uygular ve paketler için bir süzgeç süreci bulundurur.

Paketleri Güncel Tutun: Tüm paketleri düzenli olarak güncelleyin, böylece herhangi bir güvenlik açığının kapatıldığından emin olun.

Sanal Ortamları Kullanın: Paketleri ve kütüphaneleri sistemden izole etmek için sanal ortamları kullanın. Bu, kötü amaçlı paketleri daha kolay tespit edip kaldırmanıza yardımcı olur.

Kod İncelemeleri Gerçekleştirin: Bir projede kullanılan tüm paketler ve kütüphaneler üzerinde düzenli olarak kod incelemeleri yaparak kötü amaçlı kodları tespit edin.

Paket Doğrulama Araçları Kullanın: PEP 476 ve Güvenli Paket Yükleme gibi araçları kullanarak paketlerin gerçekliğini ve bütünlüğünü kurulumdan önce doğrulayın.

Geliştiricilere (Developers) Eğitim Verin: Geliştiricilere Yapay Zeka Tedarik Zinciri Saldırıları ile ilişkilendirilen riskleri ve kurulumdan önce paketleri doğrulamanın önemini anlatın.

Risk Faktörleri

Tehdit Aracıları/ Saldırı Vektörleri	Güvenlik Zafiyetleri	Etki
İstismar Edilebilirlik: 5 (Kolay) ML Uygulamasına Özel: 5 ML İşlemlere Özel: 3	Tespit Edilebilirlik: 5 (Kolay)	Teknik: 4 (Orta)
Tehdit Aktörü: Kötü niyetli saldırgan. Saldırı Vektörü: Makine öğrenimi projesi tarafından kullanılan açık kaynaklı paketin kodunu değiştirme.	Şüpheli üçüncü taraf koduna güvenme.	Makine öğrenimi projesinin tehlikeye girmesi ve organizasyona zarar oluşturma ihtimali

Bu tablonun yalnızca [aşağıdaki senaryoya](#) dayalı bir örnek olduğunu unutmamak önemlidir. Gerçek risk değerlendirmesi, her bir makine öğrenimi sisteminin kendine özgü koşullarına bağlı olacaktır.



Saldırı Senaryosu Örneği

Senaryo #1: Bir organizasyondaki makine öğrenimi projesine saldırı

Kötü niyetli bir saldırgan, büyük bir organizasyon tarafından geliştirilen bir makine öğrenimi projesini tehlikeye atmak ister. Saldırgan, projenin bir dizi açık kaynak paket ve kütüphane üzerine kurulu olduğunu biliyor ve projeyi tehlikeye atmanın bir yolunu bulmak istiyordur.

Saldırgan, bu saldırıyı projenin güvendiği paketlerden birinin kodunu değiştirerek gerçekleştirir, örneğin NumPy veya Scikit-learn gibi. Saldırgan daha sonra bu değiştirilmiş sürümü PyPI gibi genel bir depoya yükler, böylece başkalarının indirip kullanmasına olanak tanır. Saldırıya uğrayan organizasyon paketi indirip kurduğunda, saldırganın kötü amaçlı kodu da kurulur ve projeyi tehlikeye atmak için kullanılabilir.

Bu tür bir saldırı, saldırıya uğrayan organizasyonun kullandığı paketin tehlikede olduğunu fark etmediği için uzun bir süre gözden kaçabilir, bu nedenle özellikle tehlikeli olabilir. Saldırganın kötü amaçlı kodu, hassas bilgileri çalmak, sonuçları değiştirmek hatta makine öğrenme modelini başarısız kılmak için dahi kullanılabilir.

Referanslar



ML07:2023 Transfer Öğrenme Saldırısı

Tanım

Transfer öğrenme saldırıları, bir saldırganın bir modeli bir görevde eğitmesi ve ardından başka bir görevde ince ayar yaparak istenmeyen bir şekilde davranmasına neden olmaya çalıştığı durumlarda meydana gelir.

Nasıl Önlenir

Düzenli olarak eğitim veri setlerini izleyin ve güncelleyin: Eğitim veri setlerini düzenli olarak izlemek ve güncellemek, saldırganın modelinden hedef modeline kötü niyetli bilgi transferini önlemeye yardımcı olabilir.

Güvenli ve güvenilir eğitim veri setleri kullanın: Güvenli ve güvenilir eğitim veri setleri kullanmak, saldırganın modelinden hedef modeline kötü niyetli bilgi transferini önlemeye yardımcı olabilir.

Model izolasyonu uygulayın: Model izolasyonu uygulamak, bir modelden diğerine kötü niyetli bilgi transferini önlemeye yardımcı olabilir. Örneğin, eğitim ve dağıtım ortamlarını ayırarak, saldırganların eğitim ortamından dağıtım ortamına bilgi transferini engelleyebilirsiniz.

Diferansiyel gizlilik kullanın: Diferansiyel gizlilik kullanmak, eğitim veri setindeki bireysel kayıtların gizliliğini korumaya yardımcı olabilir ve saldırganın modelinden hedef modeline kötü niyetli bilgi transferini önleyebilir.

Düzenli güvenlik denetimleri yapın: Düzenli güvenlik denetimleri, sistemdeki zayıflıkları belirleyerek ve ele alarak transfer öğrenme saldırılarını tanımlamaya ve önlemeye yardımcı olabilir.

Risk Faktörleri

Tehdit Aracıları/ Saldırı Vektörleri	Güvenlik Zafiyetleri	Etki
<p>İstismar Edilebilirlik: 5 (Kolay) ML Uygulamasına Özel: 4 Saldırı özellikle makine öğrenimi uygulamasını hedef alır, dolayısıyla modele ve organizasyona zarar oluşturur. ML İşlemlere Özel: 3</p>	<p>Tespit Edilebilirlik: 1 (Zor) Saldırının tespit edilmesi tehlikede olan model tarafından üretilen sonuçların doğru gibi görünmesi ve beklentilerle uyumlu olması nedenleriyle zor olabilir.</p>	<p>Teknik: 5 (Zor) Bu saldırı, makine öğrenme alanında yüksek düzeyde teknik bilgi gerektirir ve eğitim veri setinin veya önceden eğitilmiş modellerin bütünlüğünün güvenliğini tehlikeye atma isteğine ihtiyaç duyar.</p>
<p>Bu saldırı, makine öğrenme işlemleri hakkında bilgi gerektirir, ancak zorlanmadan gerçekleştirilebilir. Tehdit Aktörü: Kötü niyetli aktör. Saldırı Vektörü: Makine öğrenme bilgisine sahip ve eğitim veri setine veya önceden eğitilmiş modellere erişimi olan saldırgan.</p>	<p>Eğitim veri seti ve önceden eğitilmiş modeller için yeterli veri koruma önlemlerinin eksikliği. Önceden eğitilmiş modellerin güvensiz saklanması ve paylaşılması. Önceden eğitilmiş modeller ve eğitim veri seti için uygun veri koruma önlemlerinin eksikliği.</p>	<p>Makine öğrenme modelinden yanıltıcı veya yanlış sonuçlar. Eğitim veri setindeki hassas bilgilerin gizliliğinin ihlali. Organizasyona itibar zararı. Hukuki sorunlar veya mevzuata uyumluluk sorunları.</p>

Bu tablonun yalnızca [aşağıdaki senaryoya](#) dayalı bir örnek olduğunu unutmamak önemlidir. Gerçek risk değerlendirmesi, her bir makine öğrenimi sisteminin kendine özgü koşullarına bağlı olacaktır.



Saldırı Senaryosu Örneği

Senaryo #1: Bir modelin kötü niyetli bir veri seti üzerinde eğitilmesi

Bir saldırgan, yüzlerin manipüle edilmiş görüntülerini içeren kötü niyetli bir veri kümesi üzerinde bir makine öğrenimi modeli eğitir. Saldırgan, bir güvenlik firması tarafından kimlik doğrulama amacıyla kullanılan yüz tanıma sistemini hedeflemek istemektedir.

Saldırgan, modelin bilgisini hedef yüz tanıma sistemine aktarır. Hedef sistem, saldırganın manipüle ettiği modeli kimlik doğrulama için kullanmaya başlar.

Sonuç olarak, yüz tanıma sistemi yanlış tahminler yapmaya başlar, bu da saldırganın güvenliği aşmasına ve hassas bilgilere erişmesine olanak tanır. Örneğin, saldırgan kendilerinin manipüle edilmiş bir görüntüsünü kullanabilir ve sistem onları meşru bir kullanıcı olarak tanır.

Referanslar



ML08:2023 Model Eğriliği

Tanım

Model eğriliği (model skewing) saldırıları, bir saldırganın eğitim verilerinin dağılımını manipüle ederek modelin istenmeyen bir şekilde davranmasına neden olduğu saldırılarla ortaya çıkar.

Nasıl Önlenir

Sağlam erişim kontrolleri uygulayın: Yalnızca yetkilendirilmiş personelin MLOps sistemi ve geri besleme döngülerine erişimi olduğundan, tüm faaliyetlerin kaydedildiğinden ve denetlendiğinden emin olun.

Geri bildirim verilerinin gerçekliğini doğrulayın: Sistem tarafından alınan geri bildirim verilerinin gerçek olduğunu doğrulamak için dijital imza ve sağlama toplamı (checksums) gibi teknikleri kullanın ve beklenen formata uymayan verileri reddedin.

Veri doğrulama ve temizleme tekniklerini kullanın: Yanlış veya kötü niyetli verilerin kullanılma riskini en aza indirme doğrultusunda geri bildirim verileriyle eğitim verilerini güncellemeden önce geri bildirim verilerini temizleyin ve doğrulayın.

Anomali tespiti uygulayın: Geri bildirim verilerindeki anormallikleri tespit etmek için istatistiksel ve makine öğrenimi tabanlı yöntemler gibi teknikleri kullanarak geri bildirim verilerindeki sapmaları tespit edin ve geri bildirim verilerindeki bu anormallik bir saldırıyı işaret edebileceğinden dolayı alarm (alert) oluşturun.

Modelin performansını düzenli olarak izleyin: Sürekli olarak modelin performansını izleyin ve tahminlerini gerçek sonuçlarla karşılaştırarak herhangi bir sapma veya eğrilik sorununun oluşup oluşmadığını kontrol edin.

Modeli sürekli olarak eğitin: Modeli güncellenmiş ve doğrulanmış eğitim verileri kullanarak düzenli olarak yeniden eğitin, böylece en güncel bilgi ve trendleri yansıtmaya devam edin.



Risk Faktörleri

Tehdit Aracıları/ Saldırı Vektörleri	Güvenlik Zafiyetleri	Etki
<p>İstismar Edilebilirlik: 5 (Kolay) ML Uygulamasına Özel: 4 Saldırgan, makine öğrenimi projesini ve projenin zayıflıklarını net bir şekilde anlamıştır.</p> <p>ML İşlemlere Özel: 3 Eğitim verilerinin manipülasyonu, makine öğrenimi sürecini bilmeyi gerektirir.</p> <p>Tehdit Aktörleri: Kötü niyetli aktörler veya bir modelin sonuçlarını manipüle etmekte çıkarı olan üçüncü taraflar.</p>	<p>Tespit Edilebilirlik: 2 (Zor) Model eğriliği, test aşamasında kolayca fark edilmeyebilir.</p> <p>Modelin, eğitim verisinin temel dağılımını doğru bir şekilde yansıtamama durumu. Bu, veri sapması(bias), verinin yanlış örnekleme yapılması veya bir saldırganın veriyi veya eğitim sürecini manipüle etmesi gibi faktörlerden kaynaklanabilir.</p>	<p>Teknik: 5 (Zor)</p> <p>Modelin çıktısına dayalı olarak yanlış kararların alınmasına yol açabilen önemli bir risktir. Bu, finansal kayıplara, itibar kaybına ve hatta model tıbbi teşhis veya cezai adalet gibi kritik uygulamalar için kullanılıyorsa bireylere zarar verme riskine neden olabilir.</p>

Bu tablonun yalnızca [aşağıdaki senaryoya](#) dayalı bir örnek olduğunu unutmamak önemlidir. Gerçek risk değerlendirmesi, her bir makine öğrenimi sisteminin kendine özgü koşullarına bağlı olacaktır.

Saldırı Senaryosu Örneği

Senaryo #1: Model eğriliği ile finansal kazanç

Bir finansal kurum, kredi başvuru sahiplerinin kredibilitelerini tahmin etmek için bir makine öğrenimi modeli kullanıyor ve modelin tahminleri kredi onay süreçlerine entegre ediyor. Bu senaryoda saldırgan, kredisinin onaylanma şansını artırmak istiyor, bu nedenle MLOps sisteminde geri bildirim döngüsünü manipüle ediyor. Saldırgan, sisteme yüksek riskli başvuru sahiplerinin geçmişte kredi onayı aldığını gösteren sahte geri bildirim verileri sunuyor ve bu geri bildirim verileri modelin eğitim verisini güncellemek için kullanılıyor. Sonuç olarak, modelin tahminleri düşük riskli başvuru sahipleri lehine eğriliyor ve saldırganın kredi onayını alma şansı önemli ölçüde artıyor.

Bu tür bir saldırı, modelin doğruluğunu ve adil olma özelliğini tehlikeye atabilir, finansal kurumun ve müşterilerinin potansiyel olarak zarar görmesine neden olabilir ve istenmeyen sonuçlara yol açabilir.

Referanslar



ML09:2023 Çıktı Bütünlüğü Saldırısı

Tanım

Bir Çıkış Bütünlüğü Saldırısı senaryosunda, bir saldırgan bir makine öğrenimi modelinin çıktısını değiştirmeyi veya manipüle etmeyi amaçlar, bu şekilde modelin davranışını değiştirmek veya kullanıldığı sisteme zarar vermek ister.

Nasıl Önlenir

Kriptografik yöntemlerin kullanılması: Dijital imza ve güvenli karma (secure hashes) gibi kriptografik yöntemler, sonuçların gerçekliğini doğrulamak için kullanılabilir.

Güvenli iletişim kanalları: Model ile sonuçları görüntülemekten sorumlu arayüz arasındaki iletişim kanalları, SSL/TLS gibi güvenli protokoller kullanılarak güvence altına alınmalıdır.

Girdi (Input) Doğrulama: Sonuçlar üzerinde beklenmeyen veya manipüle edilmiş değerleri kontrol etmek için giriş doğrulama yapılmalıdır.

Müdahaleye dayanıklı loglar: Tüm giriş ve çıkış etkileşimlerinin müdahaleye dayanıklı loglarının tutulması, çıkış bütünlük saldırılarını tespit etmeye ve buna yanıt vermeye yardımcı olabilir.

Düzenli yazılım güncellemeleri: Düzenli yazılım güncellemeleri, güvenlik açıklarını düzeltmek ve güvenlik yamalarını uygulamak için kullanılabilir ve çıkış bütünlük saldırılarının riskini azaltmaya yardımcı olabilir.

İzleme ve denetleme: Sonuçların ve model ile arayüz arasındaki etkileşimlerin düzenli izlenmesi ve denetlenmesi, şüpheli faaliyetleri tespit etmeye ve buna uygun şekilde yanıt vermeye yardımcı olabilir.

Risk Faktörleri

Tehdit Aracıları/ Saldırı Vektörleri	Güvenlik Zaafiyetleri	Etki
<p>İstismar Edilebilirlik: 5 (Kolay) ML Uygulamasına Özel: 4 ML İşlemlere Özel: 4</p> <p>Tehdit Aktörleri: Kötü niyetli saldırganlar veya modelin girdi ve çıktılarına erişimi olan iç çalışanlar ya da kişiler (insiders). Girdi ve çıktılarına erişimi olan ve belirli bir sonuca ulaşmak için bunları manipüle edebilecek üçüncü taraflar.</p>	<p>Tespit Edilebilirlik: 3 (Orta)</p> <p>Girdi ve çıktıların bütünlüğünü sağlamak için yeterli kimlik doğrulama ve yetkilendirme önlemlerinin eksikliği.</p> <p>Girdi ve çıktıların güvenliğinin genel yeterliliği hakkında yapılan performans doğrulama, standartlara uygunluğunu ölçme ve benzeri doğrulama işlemlerinin manipülasyonu önlemede yetersiz kalması.</p> <p>Girdi ve çıktıların izlenmesi (monitoring) ve loglarının tutulmasının manipülasyonu önlemede yetersiz kalması.</p>	<p>Teknik: 3 (Orta)</p> <p>Modelin tahminlerine ve sonuçlarına güvenin kaybı oluşabilir. Modelin tahminleri önemli kararlar almak için kullanılıyorsa finansal kayıp veya itibar kaybı yaşanabilir. Modelin finansal dolandırıcılık tespiti veya siber güvenlik gibi kritik bir uygulamada kullanılması durumunda güvenlik riskleri doğabilir.</p>



Bu tablonun yalnızca [aşağıdaki senaryoya](#) dayalı bir örnek olduğunu unutmamak önemlidir. Gerçek risk değerlendirmesi, her bir makine öğrenimi sisteminin kendine özgü koşullarına bağlı olacaktır.

Saldırı Senaryosu Örneği

Senaryo #1: Hasta sağlık kayıtlarının değiştirilmesi

Bir saldırgan, hastanede hastalıkları teşhis etmek için kullanılan bir makine öğrenimi modelinin çıktısına erişim kazanmıştır. Saldırgan, modelin çıktısını değiştirir ve hastalar için yanlış teşhisler sunmasını sağlar. Sonuç olarak, hastalara yanlış tedaviler verilir, bu da daha fazla zarara ve potansiyel olarak ölüme yol açabilir.

Referanslar



ML10:2023 Model Zehirleme

Tanım

Model zehirleme saldırıları, bir saldırganın modelin parametrelerini manipüle ederek onun istenmeyen bir şekilde davranmasına neden olduğu durumlarda meydana gelir.

Nasıl Önlenir

Düzenleme (Regularisation): Kayıp işlevine L1 (Lasso Düzenlemesi) veya L2 (Ridge Düzenlemesi) düzenleme gibi teknikleri eklemek, aşırı uyumu (overfitting) önlemeye ve model zehirleme saldırılarının olasılığını azaltmaya yardımcı olur.

Sağlam Model Tasarımı: Sağlam mimarilere ve aktivasyon işlevlerine sahip modellerin tasarlanması, model zehirleme saldırılarının başarı oranını azaltmaya yardımcı olabilir.

Kriptografik Teknikler: Kriptografik teknikler, modelin parametrelerini ve ağırlıklarını güvence altına almak ve bu parametrelerin yetkisiz erişimini veya manipülasyonunu önlemek için kullanılabilir.

Risk Faktörleri

Tehdit Aracıları/ Saldırı Vektörleri	Güvenlik Zaafiyetleri	Etki
<p>İstismar Edilebilirlik: 5 (Kolay) ML Uygulamasına Özel: 4 ML İşlemlere Özel: 4</p> <p>Tehdit Aktörü: Derin öğrenme modellerini manipüle etmek için bilgi ve kaynaklara sahip kötü niyetli bireyler veya organizasyonlar. Derin öğrenme modelini geliştiren organizasyon içerisindeki kötü niyetli çalışanlar(insiders).</p>	<p>Tespit Edilebilirlik: 3 (Orta)</p> <p>Modelin koduna ve parametrelerine yetersiz erişim kontrolü. Güvenli kod yazma pratiğinin yetersizliği. . Modelin faaliyetlerinin yetersiz izlenmesi (monitoring) ve loglanması.</p>	<p>Teknik: 3 (Orta)</p> <p>Modelin tahminleri istenilen sonuçlara ulaşmak için manipüle edilebilir. Model içerisindeki gizli bilgiler çıkarılabilir. Modelin tahminlerine dayalı kararlar olumsuz etkilenebilir. Organizasyonun itibarı ve güvenilirliği etkilenebilir.</p>

Bu tablonun yalnızca [aşağıdaki senaryoya](#) dayalı bir örnek olduğunu unutmamak önemlidir. Gerçek risk değerlendirmesi, her bir makine öğrenimi sisteminin kendine özgü koşullarına bağlı olacaktır.



Saldırı Senaryosu Örneği

Senaryo #1: Model zehirlenme yoluyla finansal kazanç elde etme

Bir banka, çeklerin temizlenme sürecini otomatize etmek için el yazısı karakterleri tanımlayan bir makine öğrenimi modeli kullanıyor. Model, büyüklük, şekil, eğiklik ve boşluk gibi belirli parametrelere dayalı olarak karakterleri doğru bir şekilde tanımlamak üzere tasarlanmış ve el yazısı karakterlerini içeren büyük bir veri seti üzerinde eğitilmiştir.

Bir makine öğrenimi modelini zehirlenmek isteyen bir saldırgan, modelin parametrelerini eğitim veri setindeki resimleri değiştirerek veya doğrudan modelin parametrelerini değiştirerek manipüle edebilir. Bu, modelin karakterleri farklı şekilde tanımlamasına neden olabilir. Örneğin, saldırgan parametreleri değiştirerek modelin "5" karakterini "2" karakteri olarak tanımlamasını sağlayabilir, bu da yanlış miktarların işlenmesine yol açabilir.

Saldırgan, bu zafiyeti kullanarak sahte çekleri temizleme sürecine dahil edebilir ve model, manipüle edilmiş parametreler nedeniyle bunları geçerli olarak işleyebilir. Bu, bankaya ciddi finansal kayıplara yol açabilir.

Referanslar



Teşekkürler

Katkıda Bulunanlar

Bu harika insanlara teşekkürler:

Sagar Bhure

Shain Singh

Rob van der Veer

M S Nishanth

Rick M

Harold Blankenship

RiccardoBiosas

Aryan Kenchappagol

Mikołaj Kowalczyk

Nasıl katkıda bulunulur

Bu proje, [tüm katkıları](#) kabul eden özellikleri taşıyor. Her türlü katkıya açığız!

 **AISecLab.org Türkçe Çeviri Ekibi**

Mentor: Cihan Özhan

Gözde Sarmısak

Furkan Berk Koçoğlu

Şevval Ayşe Kenar



Terimler Sözlüğü (Tamamlanmamış)

0 1 2 3 4 5 6 7 8 9 A B C D E F G H I J K L M N O P Q R S T U V W X Y Z