

- W6 - Spam Classifier -

- We will talk about ML System Design. These videos will touch on the main issues that we may face when designing a complex machine learning system.
- We will see how to strategize putting together a complex ML system.
- ilerleyen konular alakası görünebilir, bunun sebebi herşının complex ML tasarımları karşılaşılabilecek farklı sorunlara değinecek olması.
- Çok matematiksel olgulardan bahsetmeyeceğiz ama yine de, bunlar bizi büyük zaman karandıracak.

• Let's begin with prioritizing how to spend our time on what to work on

- Let's say we wanna build a Spam Classifier -

- Diyelim ki elimizde ⁽¹⁾spam ve non-spam ⁽⁰⁾olarak etiketlenmiş mailer olsun. How do we build a classifier using supervised learning?

1st Decision we must make: How do we want to represent x , that is the features of email.

- **One way to choose:** Choose 100 words indicative of spam/not spam.
E.g: deal, buy, discount, Endogen, urgent, ...

Bunlardan bazıları, spam olasılığını diğerleri ise spam olmama olasılığını temsil eder.

- Bundan sonra her email'i encode edeceğim ve bundan bir input vektörü elde edeceğim. Eğer seçilen kelimeler ilgili mailde var ise buna karşılık 1 ve yok ise 0 koyarız.

$$x = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 0 \\ \vdots \\ 0 \\ 0 \end{bmatrix} \begin{matrix} \text{deal} \\ \text{buy} \\ \vdots \\ \text{discount} \\ \vdots \\ \text{endogen} \\ \text{urgent} \end{matrix} \Rightarrow x \in \mathbb{R}^{100} : \text{en az 100 kelime seçilir}$$

\Leftarrow Mail şu olsun: "Deal of the week! Buy now!"

$$\Rightarrow x_j = \begin{cases} 1 & \text{if word } j \text{ appears in email} \\ 0 & \text{otherwise} \end{cases}$$

- Bu bir şekilde manually 100 kelime seçiyorsanız gibi göründüğü ama genellikle what's most commonly done is look through a training set, and in the training set pick the most frequently occurred n words (10000 - 50000).

W6 - Spam Classifier II -

Bu dersin

- How to spend our time to make the classifier have ^(high accuracy) low error?
 - One natural inclination (meyil) is ~~go~~ collecting lots of data
 - ↳ Ama daha önceki videolarda gördük getting a lot of data will often help but not all the time.
 - Develop Sophisticated Features based on email routing information (from email header).
 - ↳ Genelde spamler farklı headerlara sahiptir, mail farklı sehillere kullandırılır. Bunlardan Features elde edensek performansı artırabiliriz.
 - Develop sophisticated features for message body
 - ↳ E.g. should "discount" and "discounts" be treated as the same word? How about "deal" and "Dealer"? Features about punctuation?
 - Develop sophisticated algorithm to detect misspellings
 - ↳ E.g. mortgage, medicine, w4tches.

Bir sonraki videoda **Error Analysis** konseptinden bahsedeceğiz. We will talk about the way where you can try to have a more systematic way and choose amongst the options of the many different things you might work on.

- W6 - Error Analysis -

- Geçtiğimiz derste machine learning problems ile karşılaştığımızda bir çok farklı algoritma geliştirme yolu olduğundan bahsettik.
- Bu videoda ise Error Analysis kavramından bahsedeceğiz. Böylece problemlerle karşılaştığımızda daha sistematik bir yol izleyebileceğiz.

If you are starting work on a machine learning problem, or building a machine learning application there is a recommended approach:

- Instead of building a very complicated system with lots of complex features and so on, building a simple algorithm as a start that you can implement quickly is better choice. Implement something simple and quick and test it on your cross-validation data.

- Then we can plot learning curves to decide if more data, more features, etc. are likely to help. Algoritmanın high bias, high variance vb. ne problemi var gördük, bunu kullanarak ne geliştirmeye karar veririz.

• The reason that this is a good approach is often, when you are just starting on a learning problem, there is really no way to tell in advance whether you need more complex features or you need more data on something else.

• We should let evidence ~~guide~~ guide our decision on where to spend our time rather than use gut feeling, which is often wrong.

• In addition to plotting learning curves one other very useful thing is Error Analysis: Manually examine the examples (in cross validation set) that your algorithm made errors on. See if you spot any systematic trend in what type of examples it is making errors on.

↳ Often by doing that is the process that inspires you to design new features.

-WG - Error Analysis II-

Error Analysis:

- $m_{cv} = 500$ examples in cross validation set.
Algorithm misclassifies 100 emails.
Manually examine the 100 errors, and categorize them based on:
 - (i) What type of email it is (might be about pharmacy, replicas, steal password)
 - (ii) What cues (features) you think would have helped the algorithm classify them correctly.
- 1. işlem ile düzeltilmiş e-mail'leri ayırdık (hatalı olanları) ve gördük ki 100 hatanın 12'si pharmacy, 4'ü fake hesap olan, 53'ü şifre çalmaya çalışan kalan 31'i de başka şeyler. O zaman demek ki şifre çalan mail'lerinde algoritma kötü çalışıyor. Bunları daha iyi ayırabilecek feature'lar ekleyebiliriz.
- Ya da belki misspellings (meditane gibi) 4 tane, unusual email routing 16 tane, unusual punctuation 32 tane o zaman misspellings'i düzeltmek veya algılamak için advanced algorithms'e gerek yok. Bunun yerine punctuation'lara bakılabilir, buna ek olarak daha fazla gözetim.

🚫 Sonuç olarak bu şekilde bir Error Analysis ile manual olarak hata yapılan datalar incelenerek buradan sonuçlar çıkarılmaya çalışılır

Numerical Evaluation -

🚫 Lastly when developing learning algorithms, one other useful tip is to make sure that you have a numerical evaluation of your learning algorithm.

- Yani algoritmayı evaluate eden ve sonucunda tek bir real number return eden bir yapı çok yararlı olur. Belki accuracy ölçer belki error. Single real number that tells you how well your learning algorithm is doing

Importance of Numerical Evaluation: Let's say we are not sure whether or not to treat same to words: discount/discounts/discounted/discounting

- "Stemming" software ile baştaki kelimelere bakılarak benzer ise aynı kabul edilir. Bunu yapmak yararlı olabilir ama zararı da verebilir.

- universe/university karışabilir.

Thus, Error analysis may not be helpful for deciding if this is likely to improve performance. Only solution is to try and see if it works.

- Jcv ile karar verilebilir.
- Aynı şekilde uppercase-lowercase ayrılmamış mı aynı mı kabul etmeliyiz?
- Aynı şekilde veya başka bir evaluation ile karar verilebilir.
- Yine Jcv ile

- W6 - Error Metrics for Skewed Classes

For error metrics there is 1 important case, where it's particularly tricky to come up with an error/evaluation metric for the training algorithm. This case is the case called skewed (egil) classes. Let's talk about it:

Consider a Cancer Classification Example. Elimade medical patients feature olsun, bunlara göre hastalar cancer mı değil mi onu söyleriz. Malignant - benign tumor örneği gibi.

- Train a logistic regression model $h(x)$. ($y=1$ if cancer, $y=0$ o/w)
- Let's say we've test our classifier on a test set and find that we get 1% error on test (99% correct diagnoses)

★ Only 0.5% of patients training and test sets actually have cancer.

⊗ Artık 1% error iyi görünmüyor çünkü benim hipotezim sürekli 0 tahmin etse yani her gelene kanser yok dese 0.5% error oluyor daha iyi performans sergiliyormuş gibi görünüyor.

⊙ O zaman benim bu error metriğim bu tip durumlarda çok anlamlı sonuçlar vermiyor.

⊙ # of positive examples is much much smaller than the # negative ex.s. Bu yüzden bu durumun case of skewed classes diyoruz. (vice versa da olur)

⊙ Let's say we have a learning algorithm has 99.2% accuracy meaning 0.8% error. Then you made a change and the accuracy became 99.5%.

⊙ Is this an improvement or not? Did we just made an useful change or did we just replace our code with something predicts "0" more often?

⊗ So if you have skewed classes it becomes much harder to use just classification accuracy, because you can get very high classification accuracies on very low errors but it's not always clear that if you've improved your classifier performance, sadece daha fazla 0 tahmin ederek accuracy artırabiliyoruz.

Some olarak skewed classes için ayrı bir error metric kullanmalıyız.

One such evaluation metric is Precision/Recall

Precision / Recall

- When defining Precision/Recall usually we use convention: $y=1$ for the rare class (cancer)
- If we are making a binary classification there are 2 possible options for Actual Class: 0 or 1. Predicted class can be any.

Predicted Class

	Actual Class	
Predicted Class	1	0
	<div>1</div> <div>True Positive</div> <div>False Positive</div>	<div>0</div> <div>False Negative</div> <div>True Negative</div>

Precision: Of all the patients where we predicted $y=1$, what fraction actually has cancer?

$$\frac{\text{True Positives}}{\text{\# of predicted positives}} = \frac{\text{True Positives}}{\text{True Pos.} + \text{False Pos.}}$$

Recall: Of all patients that actually have cancer, what fraction did we correctly detect as having cancer?

$$\frac{\text{True Pos.}}{\text{\# of actual pos.}} = \frac{\text{True Pos.}}{\text{True Pos.} + \text{False Neg.}}$$

If a predictor always predict $y=0$ all the time and says no one has this rare condition (cancer) then the recall will be zero!

Trade-off b/w Precision and Recall

- Precision ile Recall arasında bir trade-off vardır. Bunu kontrol etmek istersen.
- Aynı example ile devam edelim logistic regression ile cancer classification

- $0 < h(x) < 1$
- Predict 1 if $h(x) \geq 0.5$
- Predict 0 if $h(x) < 0.5$

Suppose we want to predict $y=1$ (cancer) only if very confident \rightarrow avoids false positives

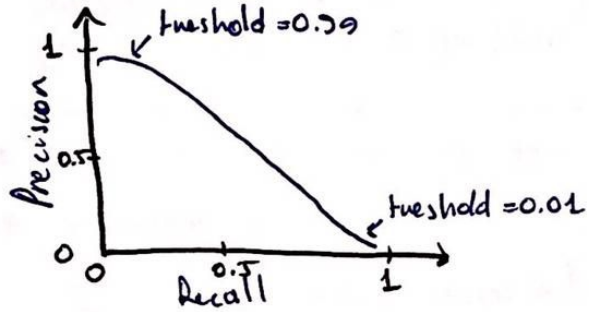
- Adama dendi ki yere kanser oldu denek ve hata yapmak istemem.
- Bunu sağlamak için bir yolu $h(x) \geq 0.9$ için $y=1$ yapmaktır. $h(x) < 0.9$ için $y=0$ predict etmektir. Böylece sadece kendine güvenince $y=1$ diyen emir olursa positive diyorum ama false negative olur çünkü 0.9 ağırlık değil diyorum.

Suppose we want to avoid missing too many cases of cancer \rightarrow avoids false negatives

- Kanseri olabilecek adama iyisin denek de cancer mal olur. İstemem.
- $h(x) \geq 0.1$ için kanser ($y=1$), $h(x) < 0.1$ için $y=0$ diyelim.
- Higher recall, lower precision.

⊗ Sonuçta $h(x) \geq \text{threshold}$ için threshold değere göre $y=1$ ve $y=0$ tahmini yapıyor isek. Bu threshold değeri Precision ile Recall trade-off'üne belirliyan.

• Precision ile Recall arasında aşağıdaki benzer bir ilişki vardır (yukarı)



• How to choose this Threshold?

• Farklı algoritmaların farklı Precision ve Recall değerleri olsun:

Alg1: $\frac{P}{0.5}$ $\frac{R}{0.4}$ which one is better?
Alg2: 0.7 0.1
Alg3: 0.02 1.0

⊗ Farklı threshold değerleri farklı Precision (P) ve Recall (R) değerleri verir. hangisinin daha iyi olduğunu karar vermek zor. Bize tek bir numerical evaluation lazım ki direkt karar verebileyim. Nasıl?

• $\frac{P+R}{2}$ olsun nasıl olur? Olmaz. Çünkü sürekli $y=1$ tahmin ederse recall = 1 olur ortalama alırken de sonucu mesela 0.51 gelir en iyisi ama aslında ise yavaş yavaş bir algoritma. Ya da her zaman $y=0$ dersen high precision low recall.

⊗ F1 Score: $2 \frac{PR}{P+R}$ if $P \geq 0$ or $R \geq 0$ $F_1 \geq 0$ yani iyi bir F_1 score için ikisi de iyi olmalı!

$F_1 = 0$ for $P=0$ and $R=0$

$F_1 = 1$ for $P=1$ and $R=1$



W6 - When to use Large Data Sets -

- We will focus on the issue of how much data to train on?
- Önceki Videolarda durduk yere gidip daha fazla training data toplamamın her zaman işe yaramayacağını görmüştük.
- Under certain conditions, getting a lot of data and training of certain type of learning algorithm can be a very effective way to get a learning algorithm to do very good performance.

- "It's not who has the best algorithm that wins. It's who has the most data."
- So when is this true and when is not?

2 conditions önemli:

1) Feature $x \in \mathbb{R}^{n+1}$ has sufficient information to predict y accurately.

- Yani kullanılan x ile gerçekleştirilen sonuç tahmin edilebilir olmalı. Bunu anlamak için sorulabilecek sorulardan biri şu: Kullanılan input ile bir human expert güvenli bir tahmin yapabilir miydi?
- Bu condition sağlanmazsa data toplamak useless.

2) Use a learning algorithm with many parameters (e.g. logistic regression / linear regression with many features; neural network with many hidden units.)

- Böylece amaç low bias, overfitting'e yatkın bir algoritma elde etmek. Overfitting de yüksek training set büyüklüğü ile elimine edilecek böylece low bias, low variance hipotezimiz olacaktır.
- Low bias algorithms

- $J_{\text{train}}(\theta)$ will be small (çünkü training seti iyi öğreniyor)

Use a very large training set (unlikely to overfit) - provides low variance

- Training set çok büyük overfittingi engeller.

- $J_{\text{train}}(\theta) \approx J_{\text{test}}(\theta)$ and $J_{\text{test}}(\theta)$ will be small!