

[SQUEAKING]

[RUSTLING]

[CLICKING]

DAVID SONTAG: OK, so then today's lecture is going to be about data set shifts, specifically how one can be robust to data set shift. Now, this is the topic that we've been alluding to throughout the semester. And the setting that I want you to be thinking about is as follows. You're a data scientist working at, let's say, Mass General Hospital, and you've been very careful in setting up your machine learning task to make sure that the data is well specified, the labels that you're trying to predict are well specified.

You train on a valid-- you train on your training data, you test it on a held-out set, you see that the model generalizes well, you do chart review to make sure what you're predicting is actually what you think you're predicting, and you even do prospective deployment where you then let your machine learning algorithm drive some clinical decision support, and you'd see things are working great.

Now what? What happens after this stage when you go to deployment? What happens when your same model is going to be used not just tomorrow but also next week, the following week, the next year?

What happens if your model, which is working well at this one hospital, then wants to-- then there's another institution, say, maybe Brigham and Women's Hospital, or maybe UCSF, or some rural hospital in the United States wants to use the same model, will it keep working in this "short term to the future" time period or in a new institution? That's the question which we're going to be talking about in today's lecture.

And we'll be talking about how one can deal with data set shift of two different varieties. The first variety is adversarial perturbations to data, and the second variety is data due to-- the data that changes for natural reasons.

Now, the reason why it's not at all obvious that your machine learning algorithm should still work in the setting is because the number one assumption we make when we do machine learning is that your training distribution, your training data, is drawn from the same distribution as your test data. So if you now go to a setting where your data distribution has changed, even if you've computed your accuracy using your held-out data and it looks good, there's no reason that should continue to look good in this new setting, where the data distribution has changed.

A simple example of what it means for a data distribution to change might be as follows. Suppose that we have as input data, and we're trying to predict some label, which maybe meant something like, why if a patient has-- or will be newly diagnosed with type 2 diabetes, and this is an example which we-- which we talked about when we introduce risk stratification, you learn a model to predict y from x .

And now suppose you go to a new institution where their definition of what type 2 diabetes means has changed. For example, maybe they don't actually have type 2 diabetes coded in their data, maybe they only have diabetes coded in their data, which is lumping together both type 1 and type 2 diabetes, type 1 being what's usually juvenile diabetes and is actually a very distinct disease from type 2 diabetes.

So now the notion of what diabetes is is different. Maybe the use case is also slightly different. And there's no reason, obviously, that your model, which was used to predict type 2 diabetes, would work for that new label.

Now, this is an example of a very type-- of a type of data set shift which is perhaps for you obvious nothing should work in the setting because here the distribution of P of y given x changes, meaning even if you have the same individual, your distribution $P(y)$ given x in, let's say, the distribution $P(0)$ and the distribution P of y given x and $P(1)$, where this is, let's say, one institution, this is another, these now are two different distributions if the meaning of the label has changed. So for the same person, there might be different distribution over what y is.

So this is one type of data shift. And a very different type of data set shift is where we assume that these two are equal. And so that would, for example, rule out this type of data set shift. But rather what changes is P of x from location 1 to location-- to location 2.

And this is the type of data set shift which will be focused on in today's lecture. It goes by the name of covariate shift. And let's look at two different examples of that.

The first example would be of an adversarial perturbation. And so we've-- you've all seen the use of convolutional neural networks for image classification problems. This is just one illustration of such an architecture.

And with such an architecture, one could then attempt to do all sorts of different object classification or image classification tasks. You could take as input this picture of a dog, which is clearly a dog. And you could modify it just a little bit. Just add in a very small amount of noise.

What I'm going to do is now I'm going to create a new image which is that original image. Now with every single pixel, I'm going to add a very small epsilon in the direction of that noise. And what you get out is this new image, which you could stare at however long you want, you're not going to be able to tell the difference. Basically to the human eye, these two look exactly identical.

Except when you take your machine learning classifier, which is trained on original unperturbed data, and now apply it to this new image, it's classified as an ostrich. And this observation was published in a paper in 2014 called "Intriguing properties of neural networks." And it really kickstarted a huge surge of interest in the machine learning community on adversarial perturbations to machine learning.

So asking questions, if you were to perturb inputs just a little bit, how does that change your classifier's output? And could that be used to attack machine learning algorithms? And how can one defend against it?

By the way, as an aside, this is actually a very old area of research. And even back in the land of linear classifiers, these questions had been studied. Although I won't get into it in this course.

So this is a type of data set shift in the sense that what we want is that this should still be classified as an ostrich-- as a dog. So the actual label hasn't changed. We would like this distribution over the labels, given the perturbed into it, to be slightly different, except that now the distribution of inputs is a little bit different because we're allowing for some noise to be added to each of the inputs.

And in this case, the noise actually isn't random, it's adversarial. And towards the end of today's lecture, I'll give you an example of how one can actually generate the adversarial image, which can change the classifier.

Now, the reason why we should care about these types of things in this course are because I expect that this type of data set shift, which is not at all natural, it's adversarial, is also going to start showing up in both computer vision and non-computer vision problems in the medical domain.

There was a nice paper by Sam Finlayson, Andy Beam, and Isaac Kohane recently, which presented several different case studies of where these problems could really arise in health care. So, for example, here what we're looking at is an image classification problem arising from dermatology. You're given as input an image.

For example, you would like that this image be classified as an individual having a particular type of skin disorder, a nevus, and this other image, melanoma. And what one can see is that with a small perturbation of the input, one can completely swap the label that would be assigned to it from one to the other. And in this paper, which we're going to post as optional readings for today's course, they talk about how one could maliciously use these algorithms for benefit.

So, for example, imagine that a health insurance company now decides in order to reimburse for an expensive biopsy of a patient's skin, a clinician or a nurse must first take a picture of the disorder and submit that picture together with the bill for the procedure. And imagine now that the insurance company were to have a machine learning algorithm be an automatic check, was this procedure actually reasonable for this condition? And if it isn't, it might be flagged.

Now, a malicious user could perturb the input such that it would, despite the patient having perhaps even completely normal-looking skin, could nonetheless be classified by a machine learning algorithm as being abnormal in some way, and thus perhaps could get reimbursed by that procedure. Now, obviously this is an example of a nefarious setting where we would then hope that such an individual would be caught by the police, sent to jail.

But nonetheless, what we would like to be able to do is build checks and balances into the system such that that couldn't even happen because to a human it's obvious that you shouldn't be able to trick-- trick anyone with such a very minor perturbation. So how do you build algorithms that could also be not tricked as easily as humans wouldn't be tricked?

AUDIENCE: Can I ask a question

DAVID Yeah.

SONTAG:

AUDIENCE: For any of these samples, did the attacker need access to the network? Is there a way to [? attack it? ?]

DAVID So the question is whether the attacker needs to know something about the function that's being used for classifying. There are examples of both what are called white box and black box attacks, where in one setting you have access to the function and other settings you don't. And so both have been studied in the literature, and there are results showing that one can attack in either setting.

Sometimes you might need to know a little bit more. Like, for example, sometimes you need to have the ability to query the function a certain number of times. So even if you don't know exactly what the function is, like you don't know the weights of the neural network, as long as you can query it sufficiently many times, you'll be able to construct adversarial examples. That would be one approach.

Another approach would be, oh, maybe we don't know the function, but we know something about the training data. So there are ways to go about doing this even if you don't perfectly know the function. Does that answer your question?

So what about a natural perturbation? So this figure just pulled from lecture 5 when we talked about non-stationarity in the context of risk stratification, that's just to remind you here the x-axis is time, that y-axis is different types of laboratory test results that might be ordered, and the color denotes how many of those laboratory tests were ordered in a certain population at a point in time.

So what we would expect to see if the data was stationary is that every row would be a homogeneous color. But instead what we see is that there are points in time, for example, a few month integrals over here, when suddenly it looks like, for some of the laboratory tests, they were never performed. That's most likely due to a data problem, or perhaps the feed of data from that laboratory test provider got lost, there were some systems problem.

But they're also going to be settings where, for example, a laboratory test is never used until it's suddenly used. And that might be because it's a new test that was just invented or approved for reimbursement at that point in time. So this is an example of non-stationarity. And, of course, this could also result in changes in your data distribution, such as what I described over there, over time.

And the third example is when you then go across institutions, wherein, of course, both the language that might be used-- you might think of a hospital in the United States versus a hospital in China, the clinical notes will be written in completely different languages, that'll would be an extreme case. And a less extreme case might be two different hospitals in Boston where the acronyms or the shorthand they use for some clinical terms might actually be different because of local practices.

So, what do we do? This is all a setup. And for the rest of the lecture, what I'll talk about is first, very briefly, how one can build in population-level checks for has something changed. And then the bulk of today's lecture, we'll be talking about how to develop transfer learning algorithms and how one could think about defenses to adversarial attacks.

So before I show you that first slide for bullet one, I want to have a bit of discussion. You've suddenly done that thing of learning machine learning algorithm in your institution, and you want to know, will this algorithm work at some other institution? You pick up the phone, you call up your collaborating data scientists at another institution, what are the questions that you should ask them when we're trying to understand, will your algorithm work there as well? Yeah.

AUDIENCE: What kind of lab test information they collect [INAUDIBLE].

DAVID So what type of data do they have on their patients, and do they have similar data types or features available for their patient population? Other ideas, someone who hasn't spoken in the last two lectures, maybe someone in the far back there, people who have their computer out. Maybe you with your hand in your mouth right there, yeah, you with your glasses on. Ideas.

[STUDENT LAUGHS]

AUDIENCE: Sorry, can you repeat the question?

DAVID You want me to repeat the question? The question was as follows. You learn your machine learning algorithm at some institution, and you want to apply it now in a new institution. What questions should you ask of that new institution to try to assess whether your algorithm will generalize in that new institution?

AUDIENCE: I guess it depends on your problem you're looking at, like whether you're trying to learn possible differences in your population, if you're requiring data with particular [INAUDIBLE] use. So I'd envision it that you'd want to, like are your machines calibrated [INAUDIBLE]? Do they use techniques to acquire the data?

DAVID All right. So let's break down each of the answers that you gave. The first answer that you gave was, are there differences in the population? What would be an example-- someone else now, what are we an example of a difference in a population? Yep.

AUDIENCE: Age distribution You might have younger people in maybe Boston versus like a Massachusetts [INAUDIBLE].

DAVID So you might have younger people in Boston versus older people who are in Central Massachusetts. How might a change in age distribution affect your ability of your algorithms to generalize? Yep.

AUDIENCE: [? Possibly ?] health patterns, where young people are very different from [INAUDIBLE] who have some diseases that are clearly more prevalent in populations that are older [? than you. ?]

DAVID Thank you. So sometimes we might expect a different just set of diseases to occur for a younger population versus an older population. So I type 2 diabetes, hypertension, these are diseases that are often diagnosed when patients-- when individuals are 40s, 50s, and older. If you have people who are in their 20s, you don't typically see those diseases in a younger population.

And so what that means is if your model, for example, was trained on a population of very young individuals, then it might not be able to-- and suppose you're doing something like predicting future cost, so something which is not directly tied to the disease itself, the features that are predictive of future cost in a very young population might be very different from features-- for predictors of cost in a much older population because of the differences in conditions that those individuals have.

Now the second answer that was given had to do with calibration of instruments. Can you elaborate a bit about that?

AUDIENCE: Yeah. So I was thinking [? clearly ?] in the colonoscopy space. But if you're collecting-- so in that space, you're collecting videos of colons. And so you can have machines that are calibrated very differently, let's say different light exposure, different camera settings. But you also have that the GIs and physicians have different techniques as to how they explore the colon. So the video data itself is going to be very different.

DAVID So the example that was given was of colonoscopies and data that might be collected as part of that. And the data that could be-- the data that could be collected could be different for two different reasons. One, because the-- because the actual instruments that are collecting the data, for example, imaging data, might be calibrated a little bit differently.

And a second reason might be because the procedures that are used to perform that diagnostic test might be different in each institution. Each one will result in slightly different biases to the data, and it's not clear that an algorithm trained on one type of procedure or one type of instrument would generalize to another. So these are all great examples.

And so when one reads a paper from the clinical community on developing a new risk stratification tool, what you will always see in this paper is what's known as "Table 1." Table 1 looks a little bit like this. Here I pulled one of my own papers that was published in *JAMA Cardiology* for 2016 where we looked at how to try to find patients with heart failure who are hospitalized. And I'm just going to walk through what this table is.

So this table is describing the population that was used in the study. At the very top, it says these are characteristics of 47,000 hospitalized patients. Then what we've done is, using our domain knowledge, we know that this is a heart failure population, and we know that there are a number of different axes that differentiate patients who are hospitalized that have heart failure. And so we enumerate over many of the features that we think are critical to characterizing the population, and we give descriptive statistics on each one of those features.

You always start with things like age, gender, and race. And so here, for example, the average age was 61 years old, this was, by the way, NYU Medical School, 50.8% female, 11.2% Black, African-American, 17.6% of individuals were on Medicaid, which was a state-provided health insurance for either disabled or lower-income individuals.

And then we looked at quantities like what types of medications were patients on. 41% of-- 42% of inpatient patients were on something called beta blockers. 31.6% of outpatients were on beta blockers.

We then looked at things like laboratory test results. So one can look at the average creatinine values, the average sodium values of this patient population. And in this way, it described what is the population that's being studied.

Then when you go to the new institution, that new institution receives not just the algorithm, but they also receive this Table 1 that describes a population in which the algorithm was learned on. And they could use that together with some domain knowledge to think through questions like what we were eliciting-- what I elicited from you in our discussion so that we could think, is it actually-- does it make sense that this model will generalize to this new institution? Are the reasons why it might not? And you could do that even before doing any prospective evaluation on the new population.

So almost all of you should have something like Table 1 in your project write-ups because that's an important part of any study in this field is describing, what is the population that you're doing your study on? You agree with me, Pete?

PETER SZOLOVITS: Yeah. I would just add that Table 1, if you're doing a case control study, you will have two columns that show the distributions in the two populations, and then a p-value of how likely those differences are to be significant. And if you leave that out, you can't get your paper published.

DAVID I'll just repeat Pete's answer for the recording. If you are-- this table is for a predictive problem. But if you're thinking about a causal inference type problem, where there's a notion of different intervention groups, then you'd be expected to report the same sorts of things, but for both the case population, the people who received, let's say, treatment one, and the control population of people who receive treatment zero. And then you would be looking at differences between those populations as well at the individual feature level as part of the descriptive statistics for that study.

Now, this-- yeah.

AUDIENCE: Is this to identify [? individually ?] [? between ?] those peoples? [INAUDIBLE] institutions to do like t-tests on those tables--

DAVID To see if they're different? No, so they're always going to be different. You go to a new institution, it's always going to look different. And so just looking to see how something changed is not-- the answer's always going to be yes.

But it enables a conversation to think through, OK, this, and then you might look-- you might use some of the techniques that Pete's going to talk about next week on interpretability to understand, well, what is the model actually using. Then you might ask, oh, OK, well, the model is using this thing, which makes sense in this population but might not make sense in another population. And it's these two things together that make the conversation.

Now, this question has really come to the forefront in recent years in close connection to the topic that Pete discussed last week on fairness in machine learning. Because you might ask if a classifier is built in some population, is it going to generalize to another population if that population that has learned on was very biased, for example, it might have been all white people. You might ask, is that classifier going to work well in another population that might perhaps include people of different ethnicities?

And so that has led to a concept which was recently published. This working draft that I'm showing the abstract from was just a few weeks ago called "Datasheets for data sets." And the goal here is to standardize the process of describing-- of eliciting the information about what is it about the data set that really played into your model?

And so I'm going to walk you through very briefly just through a couple of elements of what an example data set for a datasheet might look like. This is too small for you to read, but I'll blow up one section in just a second. So this is a datasheet for a data set called Studying Face Recognition in an Unconstrained Environment. So it's for computer vision problem.

There are going to be a number of questionnaires, which this paper that I point you to outlines. And you as the model developer go through that questionnaire and fill out the answers to it, so including things about motivation for the data set creation composition and so on.

So in this particular instance, this data set called Labeled Faces in the Wild was created to provide images that study face recognition in an unconstrained [INAUDIBLE] settings, where image characteristics such as pose, elimination, resolution, focus cannot be controlled. So it's intended to be real-world settings.

Now, one of the most interesting sections of this report that one should release with the data set has to do with how was the data preprocessed or cleaned? So, for example, for this data set, it walks through the following process.

First, raw images were obtained from the data set, and it consisted of images and captions that were found together with that image in news articles or around the web. Then there was a face detector that was run on the data set. Here were the parameters of the face detector that were used.

And then remember, the goal here is to study face detection. And so-- so one has to know, how were the-- how were the labels determined? And how would one, for example, eliminate if there was no face in this image? And so there they described how a face was detected and how a region was determined to not be a face in the case that it wasn't. And finally, it describes how duplicates were removed.

And if you think back to the examples we had earlier in the semester from medical imaging, for example in pathology and radiology, similar data set constructions had to be done there. For example, one would go to the PAC System where radiology images are stored, one would-- one would decide which images are going to be pulled out, one would go to radiography reports to figure out how do we extract the relevant findings from that image, which would give the labels for that predictive-- for that learning task. And each step there will incur some bias and some-- which one then needs to describe carefully in order to understand what might the bias be of the learned classifier.

So I won't go into more detail on this now, but this will also be one of the suggested readings for today's course. And it's a fast read. I encourage you to go through it to get some tuition for what are questions we might want to be asking about data sets that we create.

And for the rest of this semester-- for the rest of the lecture today, I'm now going to move on to some more technical issues. So we have to do it. We're doing machine learning now. The populations might be different.

What do we do about it? Can we change the learning algorithm in order to hope that your algorithm might transfer better to a new institution? Or if we get a little bit of data from that new institution, could we use that small amount of data from the new institution or a future time point in the future to retrain our model to do well in that slightly different distribution? So that's the whole field of transfer learning.

So you have data drawn from one distribution p of x and y , and maybe we have a little bit of data drawn from a different distribution q of x,y . And under the covariate shift assumption, I'm assuming that $q(x,y)$ is equal to q of x times p of y given x , namely that the conditional distribution of y given x hasn't changed. The only thing that might have changed is your distribution over x . So that's what the covariate shift assumption would assume.

So suppose that we have some small amount of data drawn from the new distribution q . How could we then use that in order to perhaps retrain our classifier to do well for that new institution? So I'll walk through four different approaches to do so. I'll start with linear models, which are the simplest to understand, and then I'll move on to deep models.

The first approach to something that you've seen already several times in this course. We're going to think about transfer as a multi-task learning problem, where one of the tasks has much less data than the other task. So if you remember when we talked about disease progression modeling, I introduced this notion of regularizing the weight vectors so that they could be close to one another. At that time, we were talking about weight vectors predicting disease progression in different time points in the future.

We could use exactly the same idea here, where you take your classifier, your linear classifier that was trained on a really large corpus, I'm going to call that-- I'm going to call the weights of that classifier w_{old} , and then I'm going to solve a new optimization problem, which is minimizing over the weights w that minimizes some loss. So this is where your training-- your new training data come in. So I'm going to assume that the new training set D is drawn from the q distribution. And I'm going to add on a regularization that asks that w should stay close to w_{old} .

Now, if the amount of data you have-- if D , the data from that new institution, was very large, then you wouldn't need this at all because you would be able to just-- you would be able to ignore the classifier that you learned previously and just refit everything to that new institution's data. Where something like this is particularly valuable is if there was a small amount of data set shift, and you only have a very small amount of labeled data from that new institution, then this would allow you to change your weight vector just a little bit. So if this coefficient was very large, it would say that the new w can't be too far from the old w . So it'll allow you to shift things a little bit in order to do well on the small amount of data that you have.

So, for example, if there is a feature which was previously predictive, but that feature is no longer present in the new data set, so, for example, it's all identically zero, then, of course, the new weight vector-- the new weight for that feature is going to be set to 0, and that weight you can think about as being redistributed to some of the other features.

Does this makes sense? Any questions? So this is the simplest approach to transfer learning. And before you ever try anything more complicated, always try this. Uh, yep.

So the second approach is also with a linear model, but here we're no longer going to assume that the features are still useful. So there might-- when you go from-- when you go from a-- your first institution, let's say, I'm GH on the left, you learn your model, and you can apply it to some new institution, let's say, UCSF on the right, it could be that there is some really big change in the feature set such that-- such that the original features are not at all useful for the new feature set.

And a really extreme example of that might be the setting that I gave earlier when I said, your model's trained on English, and you're testing it out in Chinese. That would be an example-- if you use a bag of words model, that would be an example where your model, obviously, wouldn't generalize at all because your features are completely different.

So what would you do in that setting? What's the simplest thing that you might do? So you're taking a text classifier learned in English, and you want to apply it in a setting where that language is Chinese. What would you do?

AUDIENCE: Train on them.

DAVID Translate, you said. And there was another answer.

SONTAG:

AUDIENCE: Or try train an RN.

DAVID Train an RN to do what?

SONTAG:

AUDIENCE: To translate.

DAVID Train an RN-- oh, OK. So assume that you have some ability to do machine translation, you translate from English

SONTAG: to-- from Chinese to English. It has to be that direction because the original classifier was trained in English. And then your new function is the composition of the translation and the original function, right? And then you can imagine doing some fine tuning if you had a small amount of data.

Now, the simplest translation function might be just use a dictionary. So you look up a word, and if that word has an analogy in another language, you say, OK, this is the translation. But there are always going to be some words in your language which don't have a very good translation. And so you might imagine that the simplest approach would be to translate, but then to just drop out words that don't have a good analog and force your classifier to work with, let's say, just the shared vocabulary.

Everything we're talking about here is an example of a manually chosen decision. So we're going to manually choose a new representation for the data such that we have some amount of shared features between the source and target data sets.

So let's talk about electronic health record 1 and electronic health record 2. By the way, the slides that I'll be presenting here are from a paper published in KDD by Jan, Tristan, your instructor, Pete, and John Guttag.

So you have to go two electronic health records, electronic health record 1, electronic health record 2. How can things change? Well, it could be that the same concept in electronic health record 1 might be mapped to a different encoding, so that's like an English-to-Spanish type translation, in electronic health record 2.

Another example of a change might be to say that some concepts are removed, like maybe you have laboratory test results in electronic health record 1 but not in electronic health record 2. So that's why you see an edge to nowhere. Another change might be there might be new concepts. So the new institution might have new types of data that the old institution didn't have.

So what do you do in that setting? Well, one approach we would say, OK, we have some small amount of data from electronic health record 2. We could just train using that and throw away your original data from electronic health record 1. Now, of course, if you only had a small amount of data from the target to distribution, then that's going to be a very poor approach because you might not have enough data to actually learn a reasonable enough model.

A second obvious approach would be, OK, we're going to just train on electronic health record 1 and apply it. And for those concepts that aren't present anymore, so be it. Maybe things won't work very well.

A third approach, which we were alluding to before when we talked about translation, would be to learn a model just in the intersection of the two features. And what this work does, as they say, we're going to manually redefine the feature set in order to try to find as much common ground as possible. And this is something which really involves a lot of domain knowledge. And I'm going to be using this as a point of contrast from what I'll be talking about in 10 or 15 minutes, where I talk about how one could do this without that domain knowledge that we're going to use here.

So the setting that they looked at is one of predicting outcomes, such as in-hospital mortality or length of stay. The model which is going to be used as a bag-of-events model. So we will take a patient's longitudinal history up until the time of prediction. We'll look at different events that occurred. And this study was done using PhysioNet.

And MIMIC, for example, events are encoded with some number, like 5814 might correspond to a CVP alarm, 1046 might correspond to pain being present, 25 might correspond to the drug heparin being given and so on. So we're going to create one feature for every event which has some number-- which is encoded with some number. And we'll just say 1 if that event has occurred, 0 otherwise. So that's the representation for a patient.

Now, because when one goes through this new institution, EHR2, the way that events are encoded might be completely different. One won't be able to just use the original feature representation. And that's the English-to-Spanish example that I gave.

But instead, what one could try to do is come up with a new feature set where that feature set could be derived from each of the different data sets. So, for example, since each one of the events in MIMIC has some text description that goes with it, event one corresponds to ischemic stroke, event 2, hemorrhagic stroke, and so on, one could attempt to map-- use that English description of the feature to come up with a way to map it into a common language. In this case, the common language is the UMLS, the United Medical Language System that Pete talked about a few lectures ago.

So we're going to now say, OK, we have a much larger feature set where we've now encoded ischemic stroke as this concept, which is actually the same ischemic stroke, but also as this concept and that concept, which are more general versions of that original one. So this is just general stroke, and it could be multiple different types of strokes. And the hope is that even if in-- even if the model doesn't-- even if some of these more specific ones don't show up in the new institution's data, perhaps some of the more general concepts do show up there.

And then what you're going to do is you're going to learn your model now on this expanded translated vocabulary, and then translate it. And at the new institution, you'll also be using that same common data model. And that way one hopes to have much more overlap in your feature set.

And so to evaluate this, the authors looked at two different time points within MIMIC. One time point was when the Beth Israel Deaconess Medical Center was using electronic health record called CareView. And the second time point was when that hospital was using a different electronic health record called MetaVision. So this is an example actually of non-stationarity.

Now because of them using two different electronic health records, the encodings were different. And that's why this problem arose. And so we're going to use this approach, and we're going to then learn a linear model on top of this new encoding that I just described. And we're going to compare the results by looking at how much performance was lost due to using this new encoding, and how well we generalize from one-- from one-- from the source task to the target task.

And so here's the first question, which is, how much do we lose by using this new encoding? So as a comparison point for looking at predicting in-hospital mortality, we'll look at, what is the predictive performance if you're to just use an existing, very simple risk score called the SAPS score? And that's this red line where that y-axis here is the area under the ROC curve, and the x-axis is how much time in advance you're predicting, so the prediction gap.

So using this very simple score, SAPS get somewhere between 0.75 and 0.80, area under the ROC curve. But if you were to use all of the events data, which is much, much richer than what went into that simple SAPS score, you would get the purple curve, which is-- the purple curve, which is SAPS plus the event data, or the blue curve, which is just the events data. And you can see you can get substantially better predictive performance by using that much richer feature set.

The SAPS score has the advantage that it's easier to generalize because it's so simple, those feature elements, one could trivially translate to any new EHR, either manually or automatically, and thus it'll always be a viable route. Whereas this blue curve, although it gets better predictive performance, you have to really worry about these generalization questions. And the same story happens in both of the source task and the target task.

Now the second question to ask is, well, how much do you lose when you use the new representation of the data? And so here looking at, again, both of the two-- both EHRs, what we see first in red is the same red curvature-- is the same as the blue curvature on the previous slide. It's using SAPS plus the item IDs, so using all of the data.

And then the blue curve here, which is a bit hard to see, but it's right there, it's substantially lower. So that's what happens if you now use this new representation. And you see that you do lose something by trying to find a common vocabulary. The performance does get hit a bit.

But what's particularly interesting is when you attempt to generalize, you start to see a swap. So if we now-- so now the colors are going to be quite similar. Red here was at the very top before. So red is using the original representation of the data.

Before it was at the very top. Shown here is the training error on this institution, CareView. You see, there's so much rich information in the original feature set that it's able to do very good predictive performance.

But once you attempt to translate it, so you train on CareView, but you test on MetaVision, then the test performance shown here by this solid red line is actually the worst of all of the system. So there's a substantial drop in performance because not all of these features are present in the new EHR.

On the other hand, when the translated version, despite the fact that it's a little bit worse when evaluated on the source, it generalizes much better. And so you see a significantly better performance that's shown by this blue curve here when you use this translated vocabulary. There's a question.

AUDIENCE: So would you train with full features? So how do you apply [? with ?] them if the other [? full ?] features are-- you just [INAUDIBLE].

DAVID SONTAG: So, you assume that you have come up with a mapping from the features in both of the EHRs to this common feature vocabulary of QEs. And the way that this mapping is going to be done in this paper is based on the text of the-- of the events. So you take the text-based description of the event, and you come up with a deterministic mapping to this new UMLS-based representation. And then that's what's being used. There's no fine tuning being done in this particular example.

So I consider this to be a very naive application of transfer. The results are exactly what you would expect the results to be. And, obviously, a lot of work had to go into doing this. And there's a bit of creativity in thinking that you should use the English-based description of the features to come up with the automatic mapping, but the story ends there.

And so a question which all of you might have is, how could you try to do such an approach automatically? How could we automatically find representations-- new representations of the data that are likely to generalize from, let's say, a source distribution to a target distribution? And so to talk about that, we're going to now start thinking through representation learning-based approaches, of which deep models are particularly capable.

So the simplest approach to try to do transfer learning in the context of, let's say, deep neural networks, would be to just chop off part of the network and reuse that-- some internal representation of the data in this new location. So the picture looks a little bit like this. So the data might feed in the bottom. There might be a number of convolutional layers, some fully connected layers.

And what you decide to do is you're going to take this model that's trained in one institution, you chop it at some layer, it might be, for example, prior to the last fully connected layer, and then you're going to take that-- take the new representation of your data, now the representation of the data is what you would get out after doing some convolutions followed by a single fully connected layer, and then you're going to take your target distribution's data, which you might only have a small amount of, and you learn a simple model on top of that new representation.

So, for example, you might learn a shallow classifier using a support vector machine on top of that new representation. Or you might add in some more-- a couple more layers of a deep neural network, and then fine tune the whole thing end to end. So all of these have been tried. And in some cases, one works better than another.

And we saw already one example of this notion in this course. And that was when Adam Yala spoke in lecture 13 about breast cancer and mammography, where in his approach he said that he had tried both taking a randomly initialized classifier and comparing that to what would happen if you initialized with a well-known ImageNet-based deep neural network for the problem.

And he had a really interesting story that he gave. In his case, he had enough data that he actually didn't need to initialize using this pre-trained model from ImageNet. If he had just done a random initialization, eventually-- and this x-axis, I can't remember, it might be hours of training or epochs, I don't remember, it's time-- eventually the right initialization gets to a very similar performance. But for his particular case, if you were to do a initialization with ImageNet and then fine tune, you get there much, much quicker. And so it was for the computational reason that he found it to be useful.

But in many other applications in medical imaging, the same tricks become essential because you just don't have enough data in the new test case. And so one makes use of, for example, the filters which one learns from an ImageNet's task, which is dramatically different from the medical imaging problems, and then using those same filters together with a new top layer, set of top layers in order to fine tune it for the problem that you care about. So this would be the simplest way to try to hope for a common representation for transfer in a deep architecture.

But you might ask, how would you do the same sort of thing with temporal data, not image data, maybe data that's from language, or data from time series of health insurance claims? And for that you really want to be thinking about recurrent neural networks.

So just to remind you, recurrent neural network is a recurrent architecture where you take as input some vector. For example, if you're doing language modeling, that vector might be encoding, just a one-hot encoding of what is the word at that location. So, for example, this vector might be all zeros, except for the fourth dimension, which is a 1, denoting that this word is the word, quote, "class."

And then it's fed into a recurrent unit, which takes the previous hidden state, combined it with the current input, and gets you a new hidden state. And in this way, you read in-- you encode the full input. And then you might predict-- make a classification based on the hidden state of the last time [? step. ?] That would be a common approach.

And here would be a very simple example of a recurrent unit. Here I'm using S to denote in a state. Often you will see H used to denote the hidden state. This is a particularly simple example, where there's just a single non-linearity.

So you take your previous hidden state, you hit it with some matrix $W_{s,s}$ and you add that to the input being hit by a different matrix. You now have a combination of the input plus the previous hidden state. You apply non-linearity to that, and you get your new hidden state out. So that would be an example of a typical recurrent unit, a very simple recurrent unit.

Now, the reason why I'm going through these details is to point out that the dimension of that $W_{s,x}$ matrix is the dimension of the hidden state, so the dimension of s, by the vocabulary size if you're using a one-hot encoding of the input. So if you have a huge vocabulary, that matrix, $W_{s,x}$, is also going to be equally large. And the challenge that that presents is that it would lead to overfitting on rare words very quickly.

And so that's a problem that could be addressed by instead using a low-rank representation of that $W_{s,x}$ matrix. In particular, you could think about introducing a lower dimensional bottleneck, which in this picture I'm denoting as x_t' , which is your original x_t input, which is the one-hot encoding, multiplied by a new matrix W_e . And then your recurrent unit only takes inputs of that hidden-- of that x_t' dimension, which is k, which might be dramatically smaller than v.

And you can even think about each column of that intermediate representation, We, as a word embedding. It's a way of-- and this is something that Pete talked quite a bit about when we were thinking about natural language-- when we were talking about natural language processing. And many of you would have heard about it in the context of things like Word2Vec.

So if one wanted to take a setting, for example, one institution's data where you had a huge amount of data, learn every current neural network on that institution's data, and then generalize it to a new institution, one way of trying to do that, if you think about, what is the thing that you chop, one answer might be, all you do is you keep the word embedding.

So you might say, OK, I'm going to keep the We's, I'm going to translate it back to my new institution. But I'm going to let the recurrent unit parameters-- the recurrent parameters, for example, that Ws,s you might allow it to be relearned for each new institution. And so that might be one approach of how to use the same idea that we had from feed forward networks within a recurrent setting.

Now, all of this is very general. And what I want to do next is to instantiate it a bit in the context of health care. So since the time that Pete presented the extensions of Word2Vec such as BERT and ELMo, and I'm not going to-- I'm not going to go into them now, but you can go back to Pete's lecture from a few weeks ago to remind yourselves what those were, since the time he presented that lecture, there are actually three new papers that actually tried to apply this in the health care context, one of which was from MIT.

And so these papers all have the same sort of idea. They're going to take some data set-- and these papers all use MIMIC. They're going to take that text data, they're going to learn some word embeddings or some low-dimensional representations of all words in the vocabulary.

In this case, they're not learning a static representation for each word. Instead these BERT and ELMo approaches are going to be learning-- well, you can think of it as dynamic representations. They're going to be a function of the word and their context on the left and right-hand sides.

And then what they'll do is they'll then take those representations and attempt to use them for a completely new task. Those new tasks might be on MIMIC data. So, for example, these two tasks are classification problems on MIMIC.

But they might also be on non-MIMIC data. So these two tasks are from classification problems on clinical text that didn't even come from MIMIC at all. So it's really an example of translating what you learned from one institution to another institution.

These two data sets were super small. Actually, all of these data sets were really, really small compared to the original size of MIMIC. So there might be some hope that one could learn something that really improves generalization. And indeed, that's what plays out.

So all these tasks are looking at a concept detection task. Given a clinical note, identify the segments of text within a note that refer to, for example, a disorder, or a treatment, or something else, which you then in a second stage might normalize to the UMLS.

So what's really striking about these results is what happens when you go from the left to the right column, which I'll explain in a second, and what happens when you go top to bottom across each one of these different tasks. So the left column are the results. And these results are an F score, the results, if you were to use embeddings trained on a non-clinical data set, or said definitely, not on MIMIC but on some other more general data set.

The second column is what would happen if you trained those embedding on a clinical data set, in this case, MIMIC. And you see pretty big improvements from the general embeddings to the MIMIC-based embeddings.

What's even more striking is the improvements that happen as you get better and better embeddings. So the first row are the results if you were to use just Word2Vec embeddings. And so, for example, for the I2B2 Challenge in 2010, you get 82.65 F score using Word2Vec embeddings. And if you use a very large BERT embedding, you get 90.25 F score-- F measure, which is substantially higher. And the same findings were found time and time again across different tasks.

Now, what I find really striking about these results is that I had tried many of these things a couple of years ago, not using BERT or ELMo, but using Word2Vec, and GloVe, and fastText. And what I found is that using word embedding approaches for these problems didn't-- even if you threw that in as additional features on top of other state-of-the-art approaches to this concept extraction problem, it did not improve predictive performance above the existing state of the art.

However, in this paper, here they use the simplest possible algorithm. They used a recurrent neural network fed into a conditional random field for the purpose of classifying each word into each of these categories. And the feature represent-- the features that they used are just these embedding features.

So with just the Word2Vec embedding features, the performance is crap. You don't get anywhere close to the state of art. But with the better embeddings, they actually obtain-- actually, they improved on the state of the art for every single one of these tasks. And that is without any of the manual feature engineering which we have been using in the field for the last decade. So I find this to be extremely promising.

Now you might ask, well, that is for one problem, which is classification of concepts-- or identification of concepts. What about for a predictive problem? So a different paper also published-- what month is it now, May-- so last month in April, looked at a predicted problem of 30-day readmission prediction using discharge summaries. This also was valued on MIMIC.

And their evaluation looked at the area under the ROC curve of two different approaches. The first approach, which is using a bag-of-words model, like what you did in your homework assignment, and the second approach, which is the top row there, which is using BERT embeddings, which they call Clinical BERT.

And this, again, is something which I had tackled for quite a long time. So I worked on these types of readmission problems. And bag-of-words model is really hard to beat. In fact, did any of you beat it in your homework assignment?

If you remember, there was an extra question, which is, oh, well, maybe if we used a deep learning-based approach for this problem, maybe you could get better performance. Did anyone get better performance? No.

How many of you actually tried it? Raise your hand. OK, so one-- a couple of people who are afraid to say, but yeah. So a couple of people who tried, but not many.

But I think the reason why it's very challenging to do better with, let's say, a recurrent neural network versus a bag-of-words model is because there is-- a lot of the subtlety in understanding the text is in terms of understanding the context of the text. And that's something that using these newer embeddings is actually really good at because they can get-- they could use the context of words to better represent what each word actually means. And they see substantial improvement in performance using this approach.

What about for non-text data? So you might ask when we have health insurance claims, we have longitudinal data across time. There's no language in this. It's a time series data set.

You have ICD-9 codes at each point in time, you have maybe lab test results, medication records. And this is very similar to the market scan data that you used in your homework assignment. Could one learn embeddings for this type of data, which is also useful for transfer?

So one goal might be to say, OK, let's take every ICD-9, ICD-10 code, every medication, every laboratory test result, and embed those event types into some lower dimensional space. And so here's an example of an embedding. And you see how-- this is just a sketch, by the way-- you see how you might hope that diagnosis codes for autoimmune conditions might be all near each other in some lower dimensional space, diagnosis codes for medications that treat some conditions should be near each other, and so on. So you might hope that such structure might be discovered by an unsupervised learning algorithm that could then be used within a transfer learning approach.

And indeed, that's what we found. So I wrote a paper on this in 2015/16. And here's one of the results from that paper. So this is just a look at nearest neighbors to give you some sense of whether the embedding's actually capturing the structure of the data.

So we looked at nearest neighbors of the diagnosis ICD-9 diagnosis code 710.0, which is lupus. And what you find is that another diagnosis code, also for lupus, is the first closest result, followed by connective tissue disorder, or Sicca syndrome, which is Sjogren's disease, Raynaud's syndrome, and other autoimmune conditions. So that makes a lot of sense.

You can also go across data types, like ask, what is the nearest neighbor from this diagnosis code to laboratory tests? And since we've embedded lab tests and diagnosis codes all in the same space, you can actually get an answer to that. And what you see is that these lab tests, which by the way are exactly lab tests that are commonly used to understand progression in this autoimmune condition, are the closest neighbors. Similarly, you can ask the same question about drugs and so on.

And by the way, we have made all of these embeddings publicly available on my lab's GitHub. And since the time that I wrote this paper, there have been a number of other papers, that I give citations to at the bottom here, tackling a very similar problem. This last one also put there embeddings publicly available, and is much larger than the one that we had. So these things, I think, would also be very useful as one starts to think about how one can transfer knowledge learned on one institution to another institution where you might have much less data than that other institution.

So finally I want to return back to the question that I raised in bullet two here, where we looked at a linear model with a manually chosen representation, and ask, could we-- instead of just naively chopping your deep neural network at some layer and then fine tuning, could one have learned a representation of your data specifically for the purpose of encouraging good generalization to a new institution?

And there has been some really exciting work in this field that goes by the name of Unsupervised Domain Adaptation. So the setting that's considered here is where you have data from-- you have data from first some institution, which is x comma y . But then you want to do prediction from a new institution where all you have access to at training time is x .

So as opposed to the transfer settings that I talked about earlier, now for this new institution, you might have a ton of unlabeled data. Whereas before I was talking about having just a small amount of label data, but I never talked of the possibility of having a large amount of unlabeled data. And so you might ask, how could you use that large amount of unlabeled data from that second institution in order to learn representation that actually encourages similarities from one solution to the other? And that's exactly what these domain adversarial training approaches will do.

What they do is they add a second term to the last function. So they're going to minimize-- the intuition is you're going to minimize-- you're going to try to learn parameters that minimize your loss function evaluated on data set 1. But intuitively, you're going to ask that there also be a small distance, which I'll just note as d here, between D_1 and D_2 .

And so I'm being a little bit loose with notation here, but when I calculate distance here, I'm referring to distance in representation space. So you might imagine taking the middle layer of your deep neural network, so taking, let's say, this layer, which we're going to call the feature layer, or the representation layer, and you're going to say, I want that my data under the first institution should look very similar to the data under the second institution. So the first few layers of your deep neural network are going to attempt to equalize the two data sets so that they look similar to another, at least in x space.

And we're going to attempt to find representations of your model that get good predictive performance on the data set for which you actually have the labels and for which the induced representations, let's say, the middle layer look very similar across the two data sets. And one way to do that is just to try to predict for each-- you now get a-- for each data point, you might actually say, well, which data set did it come from, data set 1 or data set 2?

And what you want is that your model should not be able to distinguish which data set it came from. So that's what it says, gradient reverse layer you want to be able to-- you want to ensure that predicting which data set that data came from, you want to perform badly on that loss functions. It's like taking the minus of that loss. And so we're not going to go into the details of that, but I just wanted to give you a reference to that approach in the bottom.

And what I want to do is just spend one minute at the very end talking now about defenses to adversarial attacks. And conceptually this is very simple. And that's why I can actually do it in one minute.

So we talked about how one could easily modify an image in order to turn the prediction from, let's say, pig to airliner. But how could we change your learning algorithm actually to make sure that, despite the fact that you do this perturbation, you still get the right prediction out, pig?

Well, to think through that, we have to think through, how do we do machine learning? Well, a typical approach to machine learning is to learn some parameters theta minimized your empirical loss. Often we use deep neural networks, which look a little like this. And we do gradient descent where we attempt to minimize some loss surfaced, find some parameters theta have as low loss as possible.

Now, when you think about an adversarial example and where they come from, typically one finds an adversarial example in the following way. You take your same loss function, now for specific input x , and you try to find some perturbation δ to x an additive perturbation, for example, such that you increase the loss as much as possible with respect to the correct label y . And so if you've increased the loss with respect to the correct label y , intuitively then when you try to see, well, what should you predict for this new perturbed input, there's going to be a lower loss for some alternative label, which is why the prediction-- the class that's predicted actually changes.

So now one can try to find these adversarial examples using the same type of gradient-based learning algorithms that one uses for learning in the first place. But what one can do is you can use a gradient descent method now-- instead of gradient descent, gradient ascent. So you take this optimization problem for a given input x , and you try to maximize that loss for that input x with this vector δ , and you're now doing gradient ascent.

And so what types of δ should you consider? You can imagine small perturbations, for example, δ that have very small maximum values. That would be an example of an L-infinity norm. Or you could say that the sum of the perturbations across, let's say, all of the dimensions has to be small. That would be corresponding to like an L1 or an L2 norm bound on what δ should be.

So now we've got everything we need actually to think about defenses to this type of adversarial perturbation. So instead of minimizing your typical empirical loss, what we're going to do is we're going to attempt to minimize an adversarial robust loss function. What we'll do is we'll say, OK, we want to be sure that no matter what the perturbation is that one adds the input, the true label y still has low loss. So you want to find parameters theta which minimize this new quantity. So I'm saying that we should still do well even for the worst-case adversarial perturbation.

And so now this would be the following new learning objective, where we're going to minimize over theta with respect to the maximum of our δ . And you have to restrict the family that these perturbations could live in, so if that δ that you're maximizing with respect to is the empty set, you get back the original learning problem.

If you let it be, let's say, all L-infinity bounded perturbations of maximum size of 0.01, then you're saying we're going to allow for a very small amount of perturbations. And the learning algorithm is going to find parameters theta such that for every input, even with a small perturbation to it, adversarially chosen, you still get good predictive performance.

And this is now a new optimization problem that one can solve. And we've now reduced the problem of finding an adversarial robust model to a new optimization problem. And what the field has been doing in the last couple of years is coming up with new optimization approaches to try to solve those problems fast.

So, for example, this paper published an ICML in 2018 by Zico Kolter and his student-- Zico just visited MIT a few weeks ago-- what it did is it said, we're going to use a convex relaxation to the rectified linear unit, which is used in many deep neural network architectures. And what it's going to do it's then going to say, OK, we're going to think about how a small perturbation to the input would be propagated in terms of getting how much that could actually change the output. And if one could be bound at every layer by layer how much a small perturbation affects the output of that layer, then one could propagate from the very bottom all the way to the loss function of the top to try to bound how much the loss function itself changes.

And a picture of what you would expect out is as follows. On the left-hand side here, you have a data point, red and blue, and the decision boundary that's learned if you didn't do this robust learning algorithm. On the right, you have now-- you'll notice a small square around each data point. That corresponds to a maximum perturbation of some limited amount.

And now you notice how the decision boundary doesn't cross any one of those squares. And that's what would be found by this learning algorithm. Interestingly, one can look at the filters that are learned by convolutional neural network using this new learning algorithm. And you find that they're much more sparse.

And so this is a very fast moving field. Every time a new adversarial attack-- every time a new adversarial defense mechanism comes up, someone comes up with a different type of attack, which breaks it. And usually that's from one of two reasons. One, because the defense mechanism isn't provable, and so one could try to come up with a theorem which says, OK, as long as you don't perturbate more than some amount, these are the results you should expect.

The other flip of the coin is, even if you come up with some provable guarantee, there might be other types of attacks. So, for example, you might imagine a rotation to the input instead of an L-infinity bounded norm that you add to it. And so for every new type of attack model, you have to think through new defense mechanisms. And so you should expect to see some iteration in the space.

And there's a website called robust-ml.org, where many of these attacks and defenses are being published to allow for the academic community to make progress here. And with that, I'll finish today's lecture.