## Evaluating Machine Learned User Experiences

Asela Gunawardana Intelligent User Experiences

Microsoft Research

## The typical machine learning problem



### Evaluation is easy: just measure $\sum_i l_i$ on the test set.

Thank You

### Questions?

**Problem:** for real problems, we need to decide what labels  $y_i$  to look at, and what loss function  $L(\cdot, \cdot)$  to use.

But is this really a serious problem?

How hard can it be?

E.g. Netflix:  $x_i = (user_i, movie_i)$   $y_i \in \{1, 2, 3, 4, 5\}$  $L(y_i, \hat{y}_i) = (y_i - \hat{y}_i)^2$ 

## Fixing the labels and loss fixes the problem

The "Netflix problem" at NIPS is:



The user's Netflix problem is:



office

STARTREN



### Popular on Netflix ${old I}$



Top Picks for Asela 🕥

### Exciting







 $(\equiv)$  $(\mathcal{P})$ . . .

Exciting

0

### Popular on Netflix 🕥

#### NETFLIX Breaking Bad NINJA SHADOW WARRIORS SAMURAI HEADHUNTERS JESSICA JUNES e tars HOUSE of CARDS STARTACT **VEPISODES** 20 GOTHAM Journey Home SOME KINDOF BEAUTIFUL LIENS CHINA'S FORBIDDEN CIT None O \* THE COMEBACK KID \*

### Ex Machina 2015 14A 108min ELERED I 5.1



A coder at a tech company wins a week-long retreat at the compound of his company's CEO, where he's tasked with testing a new artificial intelligence.

\*\*\*\*

+ Add to My List

tarr	ing		
	Domhnall Gleeson	Alicia Vikano	der
	Sonoya Mizuno	Oscar Isaac	Claire Selby
	Symara Templeman	Gana Baya	rsaikhan

Director

Alex Garland

#### Genres

Dramas Thrillers Psychological Thrille Sci-Fi Thrillers British Movies

Sci-Fi & Fantasy Sci-Fi Dram

# Does our formulation of the problem really help users find things to watch?

# Does predicting ratings help users find things to watch?



Rating

### Predicting Ratings ≠ Predicting Usage



### Predicting Ratings ≠ Predicting Usage



### Predicting Ratings ≠ Predicting Usage



## Lesson:

The "standard," "given," or "commonly used" labels and loss functions may tell us very little about how useful the system is.

 $\equiv$  $(\mathcal{A})$ 

Exciting

0

### Popular on Netflix 🕥

### NETFLIX Breaking Bad NINJA SHADOW WARRIORS SAMURAI HEADHUNTERS JESSICA JUNES what MEN HOUSE of CARDS STARTART **VEPISODES** 20 GOTHAM Journey Home SOME KINDOF BEAUTIFUL LIENS CHINA'S FORBIDDEN CIT None 0 \* THE COMEBACK KID \*

Exciting

0

Journey Home

### Popular on Netflix 🕥

#### NETFLIX Breaking Bad NINJA SHADOW WARRIORS SAMURAI HEADHUNTERS JESSICA JUNES R leca ercion MEN HOUSE of CARDS STARTREN V EPISODES 20 GOTHAM SOME KINDOF BEAUTIFUL LIENS CHINA'S FORBIDDEN CIT None O \* THE COMEBACK KID \*

Exciting

( . . .

### Popular on Netflix 🕥



 $(\equiv)$  $(\mathcal{P})$ 

Exciting

0

. . .

### Popular on Netflix 🕥



 $(\equiv)$  $(\mathcal{P})$ . . .

Exciting

0

### Popular on Netflix 🕥

#### NETFLIX **Breaking** Bad NINJA SHADOW WARRIORS SAMURAI HEADHUNTERS JESSICA JUNES $\bigcirc$ 6 MEN HOUSE of CARDS STARTACK **VEPISODES** 20 GOTHAM Journey Home SOME KINDOF BEAUTIFUL LIENS CHINA'S FORBIDDEN CIT None O \* THE COMEBACK KID \*

- 1. Log usage (not just ratings)
- 2. Train recommender on log data from before yesterday.
- 3. Recommend items for yesterday's users.
- 4. Score against yesterday's actual usage data:

	Actually Used	Actually Unused
Recommended	True Positive	False Positive
Not Recommended	False Negative	True Negative

- 1. Log usage (not just ratings)
- 2. Train recommender on log data from before yesterday.
- 3. Recommend items for yesterday's users.
- 4. Score against yesterday's actual usage data:

	Actually Used	Actually Unused
Recommended	True Positive	False Positive
Not Recommended	False Negative	True Negative

Problem:

False Positive/True Negative

Maybe the user didn't know about the video, would have happily watched it if we actually recommended it.

- 1. Log usage (not just ratings)
- 2. Train recommender on log data from before yesterday.
- 3. Recommend items for yesterday's users.
- 4. Score against yesterday's actual usage data:

	Actually Used	Actually Unused
Recommended	True Positive	False Positive
Not Recommended	False Negative	True Negative

Problem:

False Positive/True Negative

Maybe the user didn't know about the video, would have happily watched it if we actually recommended it.

## Problem #1

Our data isn't an i.i.d. draw – it's collected from a real running system.





Exciting

**()** PB

3

Journey Home

### Popular on Netflix $\Im$







Exciting

**()** PB

### Popular on Netflix ${old D}$



- 1. Log usage (not just ratings)
- 2. Train recommender on log data from before yesterday.
- 3. Recommend items for yesterday's users.
- 4. Score against yesterday's actual usage data:

	Actually Used	Actually Unused
Recommended	True Positive	False Positive
Not Recommended	False Negative	True Negative

Problems:

False Positive/True Negative

Maybe the user didn't know about the video, would have happily watched it if we actually recommended it.

**True Positive** 

Maybe the user would have watched the video already, even if we didn't predict it.

- 1. Log usage (not just ratings)
- 2. Train recommender on log data from before yesterday.
- 3. Recommend items for yesterday's users.
- 4. Score against yesterday's actual usage data:

	Actually Used	Actually Unused
Recommended	True Positive	False Positive
Not Recommended	False Negative	True Negative

Problems:

False Positive/True Negative

Maybe the user didn't know about the video, would have happily watched it if we actually recommended it.

**True Positive** 

Maybe the user would have watched the video already, even if we didn't predict it.

## Problem #2

Measuring prediction accuracy doesn't tell us how the system will **influence user behavior**.



About 12,900,000 results (0.35 seconds)

#### Mesothelioma Diagnosis?

#### Ad www.simmonsfirm.com/Mesothelioma -

We are Committed to Helping You Get the Compensation You Deserve. Why File a Lawsuit? - Mesothelioma Settlements - Where we Serve

#### Mesothelioma Diagnosis? - sokolovelaw.com Ad www.sokolovelaw.com/mesothelioma (888) 940-5538

You Didn't Deserve This Disease. Learn About Your Legal Options Now. Services: Legal Consultation, Help For Caregivers, Help For Vetearns, Mesotheliom... "A+ Rating" – Better Business Bureau Settlement Stories - Why Sokolove Law? - Free Consultation

### Mesothelioma Diagnosis? - bergmanlegal.com

Don't Be Fooled By A Claims Center. Receive Top Compensation. Call Now. Rated WA Super Lawyers · Exclusively Mesothelioma · \$500 Mil+ in Settlements Settlement Stories - Free Claim Evaluation - Reviews ♥ 821 2nd Ave #2100, Seattle, WA - Open today · 8:00 AM – 5:30 PM ▼

### Mesothelioma Cancer | Diagnosis, Treatment and Support www.mesothelioma.com/mesothelioma/

**Mesothelioma** is an aggressive cancer affecting the membrane lining of the lungs and abdomen. Malignant **mesothelioma** is the most serious of all ... Mesothelioma Symptoms - Pleural Mesothelioma - Mesothelioma Causes - Biopsies

### Mesothelioma - Overview of Malignant Mesothelioma Cancer

www.asbestos.com/mesothelioma/ - Asbestos.com

Nov 10, 2015 - Malignant **mesothelioma** is a rare, asbestos-related cancer that forms on the thin protective tissues that cover the lungs and abdomen.... Latency Period: It can take 20-50 years after asbestos exposure for **mesothelioma** to develop.... That's because this cancer can take anywhere from 20 to ...

#### Mesothelioma - Wikipedia, the free encyclopedia

https://en.wikipedia.org/wiki/**Mesothelioma** ▼ Wikipedia ▼ **Mesothelioma** (or, more precisely, malignant **mesothelioma**) is a rare form of cancer that develops from cells of the mesothelium, the protective lining that covers ... Asbestos - Mesothelium - Peritoneal mesothelioma - Category:Mesothelioma



More about this condition

#### Ads

Mesothelioma Diagnosis? www.mesotheliomaclaimscenter.info/ (855) 279-0121 You Didn't Deserve This Disease. You May Be Entitled To Compensation

#### Mesothelioma Symptoms

www.mesotheliomaportal.com/Health ▼ Search For Mesothelioma Symptoms. Learn More at Mesothelioma Portal.

#### Vou Don't Have To Sue



#### Mesothelioma Diagnosis?

#### Ad www.simmonsfirm.com/Mesothelioma •

We are Committed to Helping You Get the Compensation You Deserve. Why File a Lawsuit? - Mesothelioma Settlements - Where we Serve

#### Mesothelioma Diagnosis? - sokolovelaw.com Ad www.sokolovelaw.com/mesothelioma (888) 940-5538

You Didn't Deserve This Disease. Learn About Your Legal Options Now. Services: Legal Consultation, Help For Caregivers, Help For Vetearns, Mesotheliom... "A+ Rating" – Better Business Bureau Settlement Stories - Why Sokolove Law? - Free Consultation

### Mesothelioma Diagnosis? - bergmanlegal.com

Don't Be Fooled By A Claims Center. Receive Top Compensation. Call Now. Rated WA Super Lawyers · Exclusively Mesothelioma · \$500 Mil+ in Settlements Settlement Stories - Free Claim Evaluation - Reviews ♥ 821 2nd Ave #2100, Seattle, WA - Open today · 8:00 AM – 5:30 PM ▼

#### Mesothelioma Cancer | Diagnosis, Treatment and Support www.mesothelioma.com/mesothelioma/

**Mesothelioma** is an aggressive cancer affecting the membrane lining of the lungs and abdomen. Malignant **mesothelioma** is the most serious of all ... Mesothelioma Symptoms - Pleural Mesothelioma - Mesothelioma Causes - Biopsies

#### Mesothelioma - Overview of Malignant Mesothelioma Cancer www.asbestos.com/mesothelioma/ Asbestos.com

Nov 10, 2015 - Malignant **mesothelioma** is a rare, asbestos-related cancer that forms on the thin protective tissues that cover the lungs and abdomen.... Latency Period: It can take 20-50 years after asbestos exposure for **mesothelioma** to develop.... That's because this cancer can take anywhere from 20 to ...

#### Mesothelioma - Wikipedia, the free encyclopedia

https://en.wikipedia.org/wiki/**Mesothelioma** ▼ Wikipedia ▼ **Mesothelioma** (or, more precisely, malignant **mesothelioma**) is a rare form of cancer that develops from cells of the mesothelium, the protective lining that covers ...

Asbestos - Mesothelium - Peritoneal mesothelioma - Category:Mesothelioma

Mesothelion	na	
ABOUT	SYMPTOMS	TREATMENT
A tumor of the tissue organs.	that lines the lungs, stoma	ach, heart, and othe
Very rare Fewer than 20,000 U	S cases per year	
📋 Requires a medical diagnosis		
Lab tests or im	aging always required	
Consult a doctor for n Sources: Mavo Clinio	nedical advice c and others.	

More about this condition

#### Ads

Mesothelioma Diagnosis? www.mesotheliomaclaimscenter.info/ (855) 279-0121 You Didn't Deserve This Disease. You May Be Entitled To Compensation

Mesothelioma Symptoms www.mesotheliomaportal.com/Health ▼ Search For Mesothelioma Symptoms. Learn More at Mesothelioma Portal.

Vou Don't Have To Sue

### Mesothelioma Portal

Mesothelioma Claims – 1 (800) 713-6692

#### www.nationalmesotheliomaclaims.com

You don't have to sue anyone for financial compensation. The \$30 Billion Asbestos Trust Fund was established to pay asbestos damage claims to mesothelioma patients and their...

📞 Call

#### Mesothelioma Cancer Book - 1-800-301-1845

#### www.mesothelioma-answer.org

This authoritative book, *101 Facts About Mesothelioma* by Anna Kaplan, M.D. has received many positive reviews for providing clear and simple answers to the most common questions about mesothelioma cancer.

#### Mesothelioma Attorney Locator - 1-800-314-2433

#### www.mesothelioma-attorney-locators.com

This site allows you to easily find **Mesothelioma** attorneys located in your state. Patient advocates also help patients and their families understand what they need to know when selecting a law firm that will...

#### Mesothelioma Doctor Match - 1 (888) 888-4051

#### www.MesotheliomaDoctorMatch.com

Now you can find a top mesothelioma doctor quickly. They also help you to get an appointment more quickly while others may need to wait weeks or months to be seen. Often this can shorten the time before your treatment starts which...

#### Mesothelioma & Veterans - 1 (800) 726-7245

#### mesothelioma.veteransupport.us

Important information for Navy and other military vets who have been exposed to asbestos and now have health problems. Learn about options for making **mesothelioma** and asbestos cancer claims. Filing for veteran's benefits that can provide...

#### Official Mesothelioma Calculator - 1-800-818-7146

#### www.mesothelioma-case-calculator.com

The amount of compensation recovered in a mesothelioma case depends on what asbestos products you were exposed to, which state the case is filed in, and what law firm you hire. Some cases settle for...

#### Mesothelioma Survival Rate - 1 (888) 888-1830

#### www.mesotheliomasurvivalrate.com

How long do mesothelioma patients live after being diagnosed? What options are their for treatment? Can changes in diet help? This site answers these questions and others related to





101 Facts About Mesothelioma by Dr. Kaplan

Mesothelioma Patients & their Families Need to Know

Click Here to Get Your FREE Book!



Search

- 1. Log usage (not just ratings)
- 2. Train recommender on log data from before yesterday.
- 3. Recommend items for yesterday's users.
- 4. Score against yesterday's actual usage data:

	Actually Used	Actually Unused
Recommended	True Positive	False Positive
Not Recommended	False Negative	True Negative

Problems:

False Positive/True Negative

Maybe the user didn't know about the video, would have happily watched it if we actually recommended it.

**True Positive** 

Maybe the user would have watched the video already, even if we didn't predict it.

False Negative

Maybe the user watched the video but hated it.

- 1. Log usage (not just ratings)
- 2. Train recommender on log data from before yesterday.
- 3. Recommend items for yesterday's users.
- 4. Score against yesterday's actual usage data:

	Actually Used	Actually Unused
Recommended	True Positive	False Positive
Not Recommended	False Negative	True Negative

Problems:

False Positive/True Negative

Maybe the user didn't know about the video, would have happily watched it if we actually recommended it.

**True Positive** 

Maybe the user would have watched the video already, even if we didn't predict it.

False Negative

Maybe the user watched the video but hated it.

## Problem #3

The influence of our system may only manifest over the long term.

- 1. Our data isn't an i.i.d. draw it needs to be collected from a real running system.
- 2. Measuring prediction accuracy doesn't tell us how the system will **influence user behavior**.
- 3. The influence of our system may only manifest over the long term.

How do we avoid being fooled about how useful our system is?

### How not to be fooled

- 1. Identify what the goal is
  - Service usage
  - Sales
  - Ad monetization
  - User retention
- 2. Randomly assign users to a control and treatment group and measure improvement due to system, over time.
- 3. Use (with care) offline experiments to prioritize which experiments to run.

## The objections

Experiments are expensive and time-consuming—can only try a handful of variations.

We can't really expect scientists to build user-facing systems before they do science.

Besides, I'm are confident that <insert loss function here> will generally track <insert real criterion here>.

RMSE was good enough for Netflix: \$1,000,000 says so.

The system owner is happy with improvements in my metric.



Science is a bit like the joke about the drunk who is looking under a lamppost for a key that he has lost on the other side of the street, because that's where the light is. It has no other choice.

> Noam Chomsky (at least, according to the web)

Science is a bit like the joke about the drunk who is looking under a lamppost for a key that he has lost on the other side of the street, because that's where the light is. **It has no other choice.** 

> Noam Chomsky (at least, according to the web)

### Another choice: Build a new lamppost

(or at least a flashlight)

Joachims, KDD 2002→WSDM 2015: Use actual user behavior and mild assumptions about it to evaluate web search ranking.

Marlin et al, IJCAI 2011: How to estimate and account for selection bias in data sets.

Bottou et al, JMLR 2013: How to use data reweighting and a priori causal knowledge to correct for selection bias and make counter-factual inferences.

These issues have started to be addressed, and we need to more work that builds on this start.

### Need data that

is collected through randomization of a real system records what was presented to the user ("impression logs") records why (inputs and sampling probability/density) records what the user did