

Welcome

- Course webpage

- <http://ttic.edu/intromlss2018>

- Daily schedule (most days)

- Lectures

- 9am-noon
 - Room 526, 6045 S Kenwood Ave
 - Instructors: Suriya Gunasekar, Karl Stratos

- Lunch

- Noon-1pm

- Invited talks + Coding assignments

- 1-5pm

- Room 526+530, 6045 S Kenwood Ave

Introduction to Machine
Learning Summer School
June 18, 2018 - June 29,
2018, Chicago



Today's schedule

- 9am-9:30pm Course setup
- 9:30:00-11:00am Lecture 1.a
 - Introduction, supervised learning
- 11:00-12:30pm Lecture 1.b
 - Linear regression
- 12:30-2:00pm Lunch
 - group delivery order, or
 - walk to nearby places - Press Café, Booth School of Business, Hutchinson Commons, Build Coffee
- 2:00-5pm Programming

Teaching Assistants



Rebecca Kotsonis



Pedro Savarese



Kevin Stangl

Emails TAs ta-intromlss2018@ttic.edu

Day 1: Introduction, Linear Regression

**Introduction to Machine Learning Summer School
June 18, 2018 - June 29, 2018, Chicago**

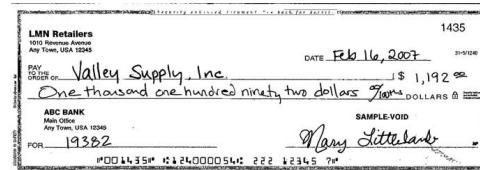
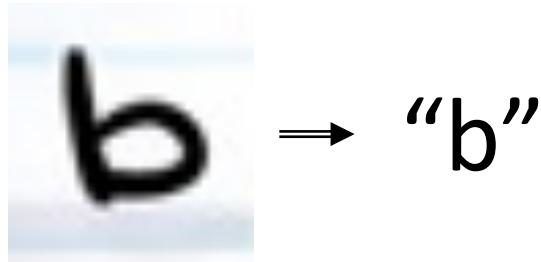
Instructor: Suriya Gunasekar, TTI Chicago

18 June 2018

What is machine learning?

Let us first see some examples...

Optical character recognition



Advantages of data driven approach

- Less programming
- Easily adaptive
- Less dependent on expert knowledge

Spam detection

Suriya Gunasekar <suriyag88@gmail.com>

Don't forget to take your free ride!Divvy Bikes <support@divvybikes.com>
Reply-To: Divvy Bikes <support@divvybikes.com>
To: suriyag88@gmail.com

Fri, May 25, 2018 at 12:24 PM

First ride FREE!



Your free ride awaits.
Download the new Divvy app to redeem.

You keep saying you're going to "explore your city more," and now we're giving you a super easy way to do it. With your first 30-minute Single Ride FREE, you'll actually hit up all those summer hot spots. Download the new Divvy app and enter code **FIRSTRIDEFREE** before this offer rides away.

[GET MY FREE RIDE](#)

Offer expires 6/17

sender info, recipients info, length of text/subject, timestamp, subject line, words/pairs of words in text, previous interactions with senders/fellow recipients



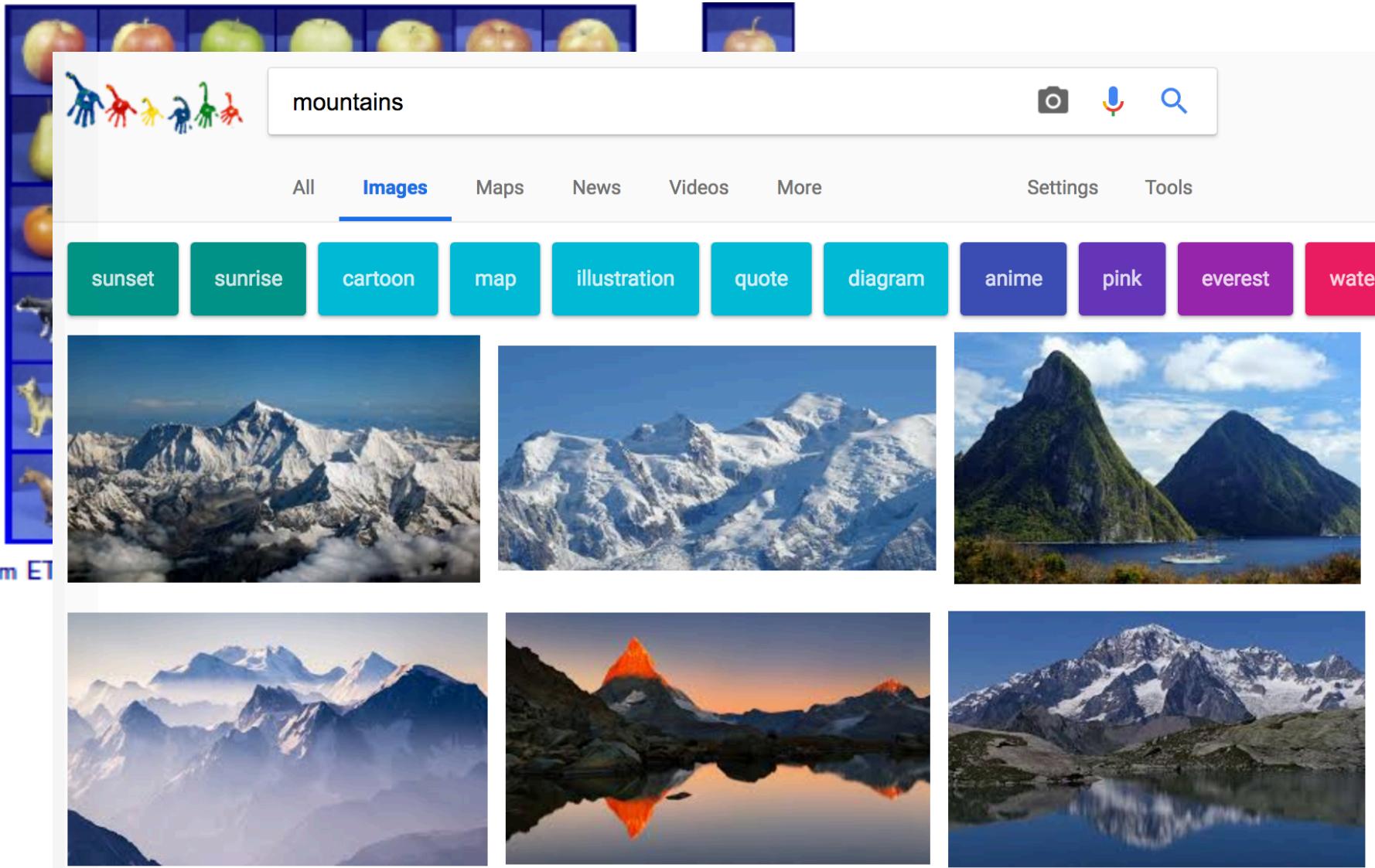
spam or
no-spam

Object recognition



From ETH database of object categories, [Leibe & Schiele 2003]

Object recognition



Face detection/recognition



image pixels



detect faces
recognize faces
detect orientation
A is talking B

Machine translation

Translate

Turn off instant translation



English Spanish French Detect language ▾



English Spanish Dutch ▾

Translate

We are learning how to use machine learning to improve our methods.



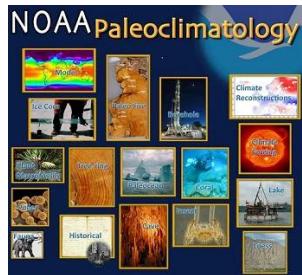
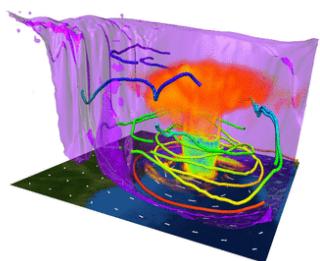
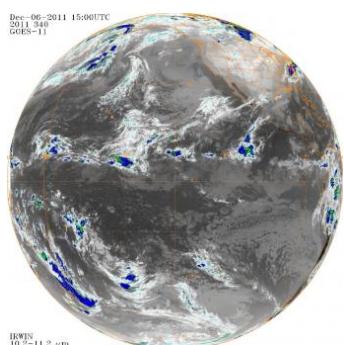
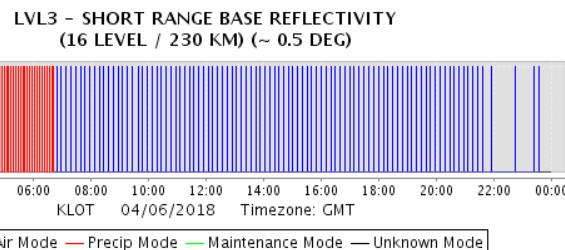
We leren hoe je machine learning te gebruiken om onze werkwijze te verbeteren.



Suggest an edit

- Complex output
- Learn from:
 - annotated translations
 - matching text
 - text (in a single language)
 - corrections

Weather forecasting



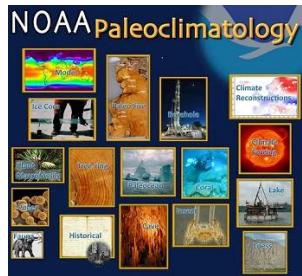
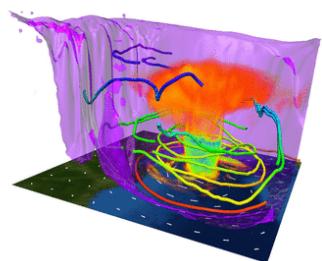
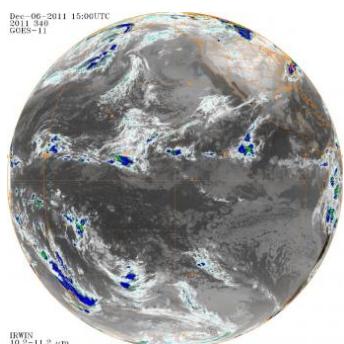
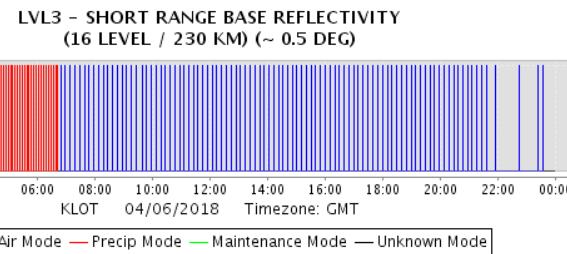
Measurements from land observatories, weather balloons, radar equipment, meteorological satellites, marine observatories, simulations, paleoclimatology data, and the historical records of these measurements



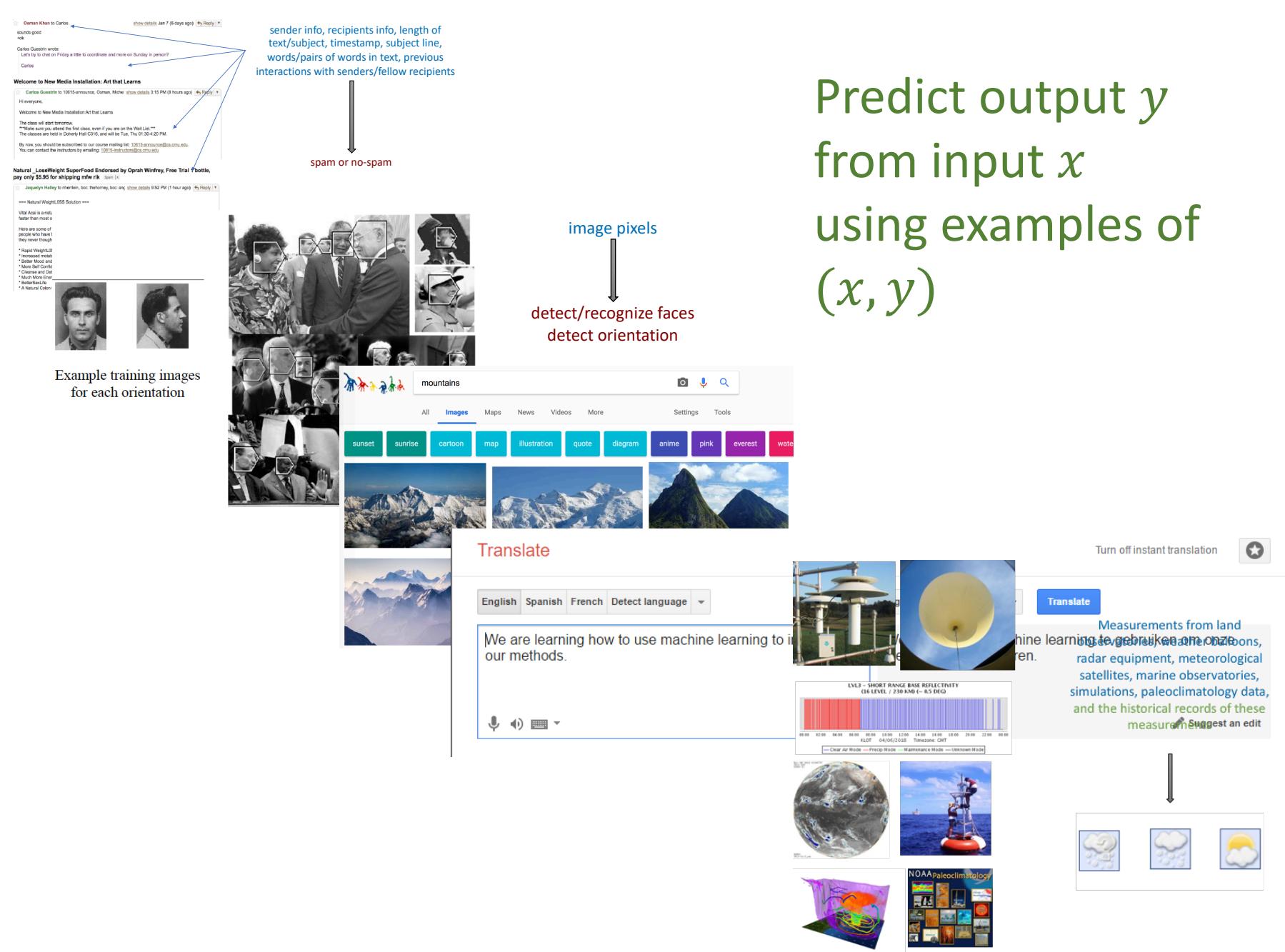
Temperature

72° F

Weather forecasting



Still largely
done by
expert build
systems





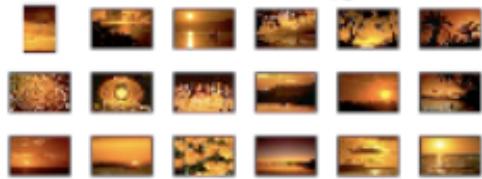
C_1



C_2



C_3

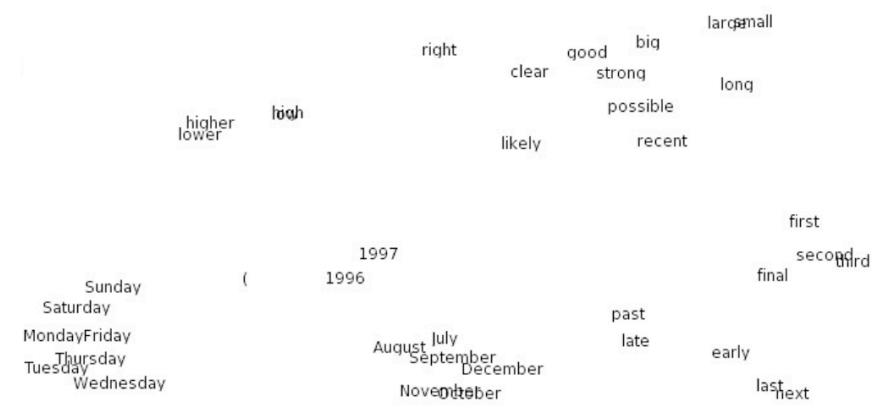


C_4

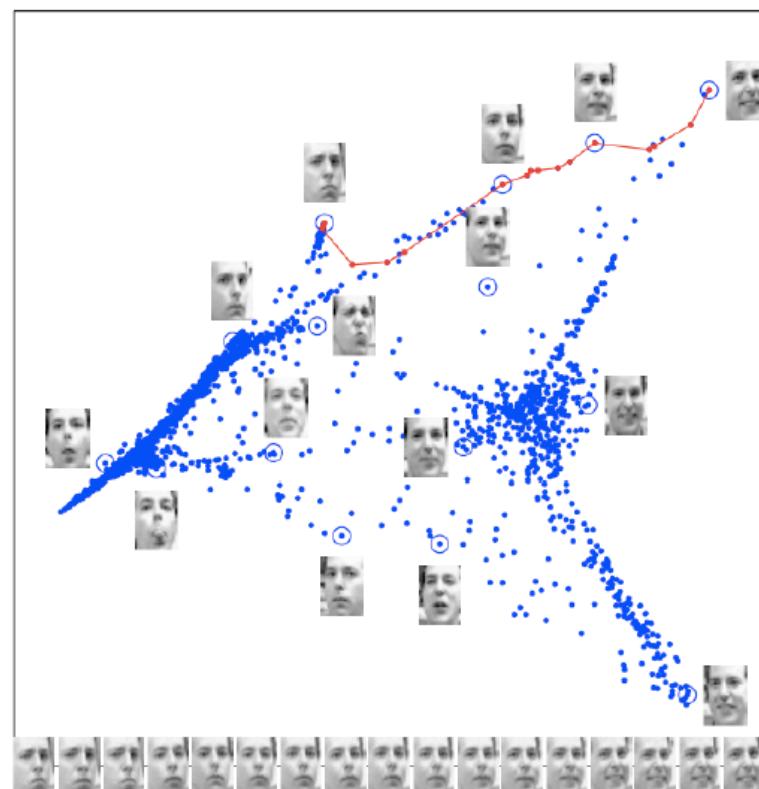


C_5

[Goldberger et al.]

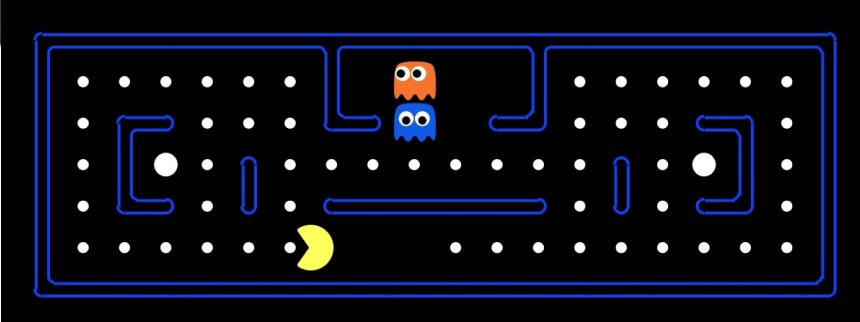


t-SNE visualization from Turian et al. (2010)



[Saul & Roweis '03]

Slide credit: David Sontag



SCORE: 18

$$\begin{array}{l}
 25: \quad 3 + 4 + p + 5 + 2p + 3 + 2p + 2p \\
 54: \quad 2 + 3 + p + 2p + 3 + 2p + 2p \\
 137: \quad p + 2 + p + 2p + 2p + 2p + 2p \\
 434: \quad 3 + 4 + p + 2p + 3 + 4 + p + 2p \\
 6130: \quad 3p + 3 + 2p + 2p + 3p + p + 2p \\
 4425: \quad 2p + p + 2p + 2p + 2p + 2p \\
 2605: \quad 2p + p + 2p + 2p + 2p + 2p \\
 3965: \quad 2p + p + 2p + 2p \\
 6573: \quad 2p + p + 2p + 2p + 2p + 2p \\
 6871: \quad 2p + p + 2p + 2p
 \end{array}$$



So what is machine learning?

What is (machine) learning?

- T. Hastie, R. Tibshirani, J. Friedman: “extract important patterns and trends, and understand ‘what the data says’”
- T. Mitchell: “The field of Machine Learning seeks to answer the question “How can we build computer systems that automatically improve with experience, and what are the fundamental laws that govern all learning processes?””
- S. Shalev-Shwartz and S. Ben-David: “automated detection of meaningful patterns in the data”
- K. Murphy: “set of methods that can automatically detect patterns in data, and then to use the uncovered patterns to predict future data or other kinds of decision making under uncertainty”
- M. Mohri, A. Rostamizadeh, and A. Talwalkar “computational methods using experience to improve performance or to make accurate predictions”

What is (machine) learning?

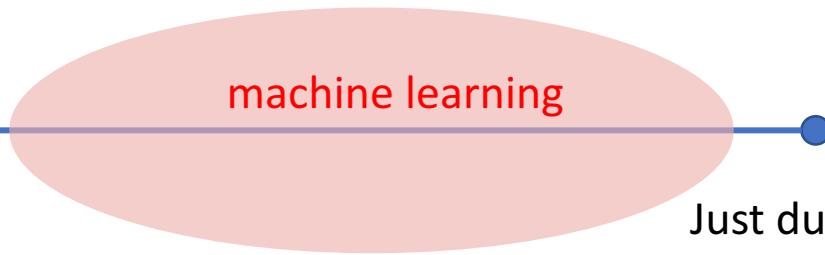
- T. Hastie, R. Tibshirani, J. Friedman: “**extract important patterns** and trends, and understand ‘what the **data** says’”
- T. Mitchell: “The field of Machine Learning seeks to answer the question “How can we build **computer systems** that **automatically** improve with **experience**, and what are the fundamental laws that govern all learning processes?””
- S. Shalev-Shwartz and S. Ben-David: “**automated detection of meaningful patterns** in the **data**”
- K. Murphy: “**set of methods** that can **automatically detect patterns** in **data**, and then to use the uncovered patterns to **predict future data** or other kinds of decision making under uncertainty”
- M. Mohri, A. Rostamizadeh, and A. Talwalkar “**computational methods** using **experience** to improve performance or to **make accurate predictions**”

Expert designed → data driven

Adder: $(x_1, x_2) \rightarrow$
 $x_1 + x_2$

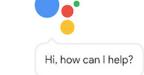
Expert
designed
systems

C. M. Bishop: “...a
training set is used to
tune the parameters of
an adaptive model”



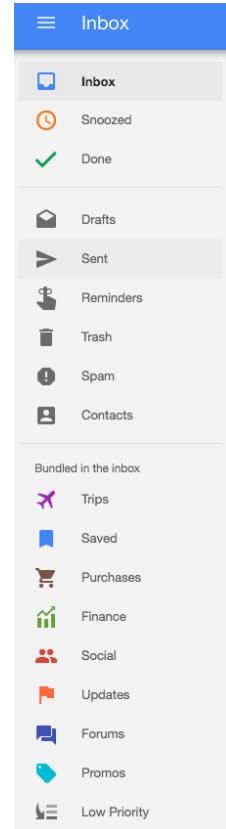
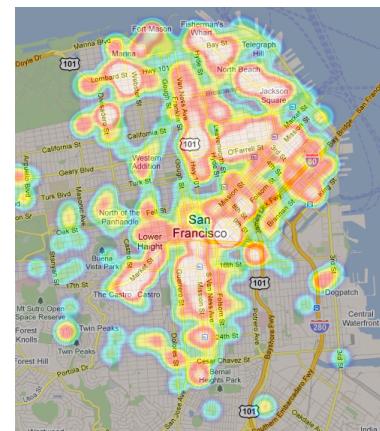
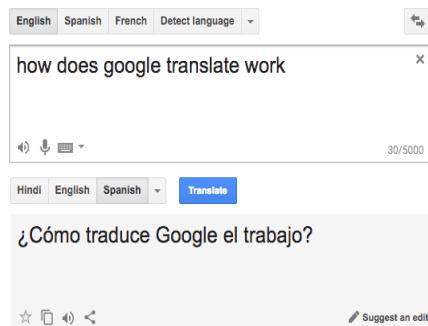
Machine learning everywhere

- Email categorization,
- Web search
- Machine translation
- Question answering, dialog systems
- Automatic speech recognition, speech synthesis
- Surveillance
- Recommendations
- Fraud protection
- Traffic routing



Meet your Google Assistant.

Ask it questions. Tell it to do things.
It's your own personal Google, always ready to help.



.....

Why machine learning?

Why automation?

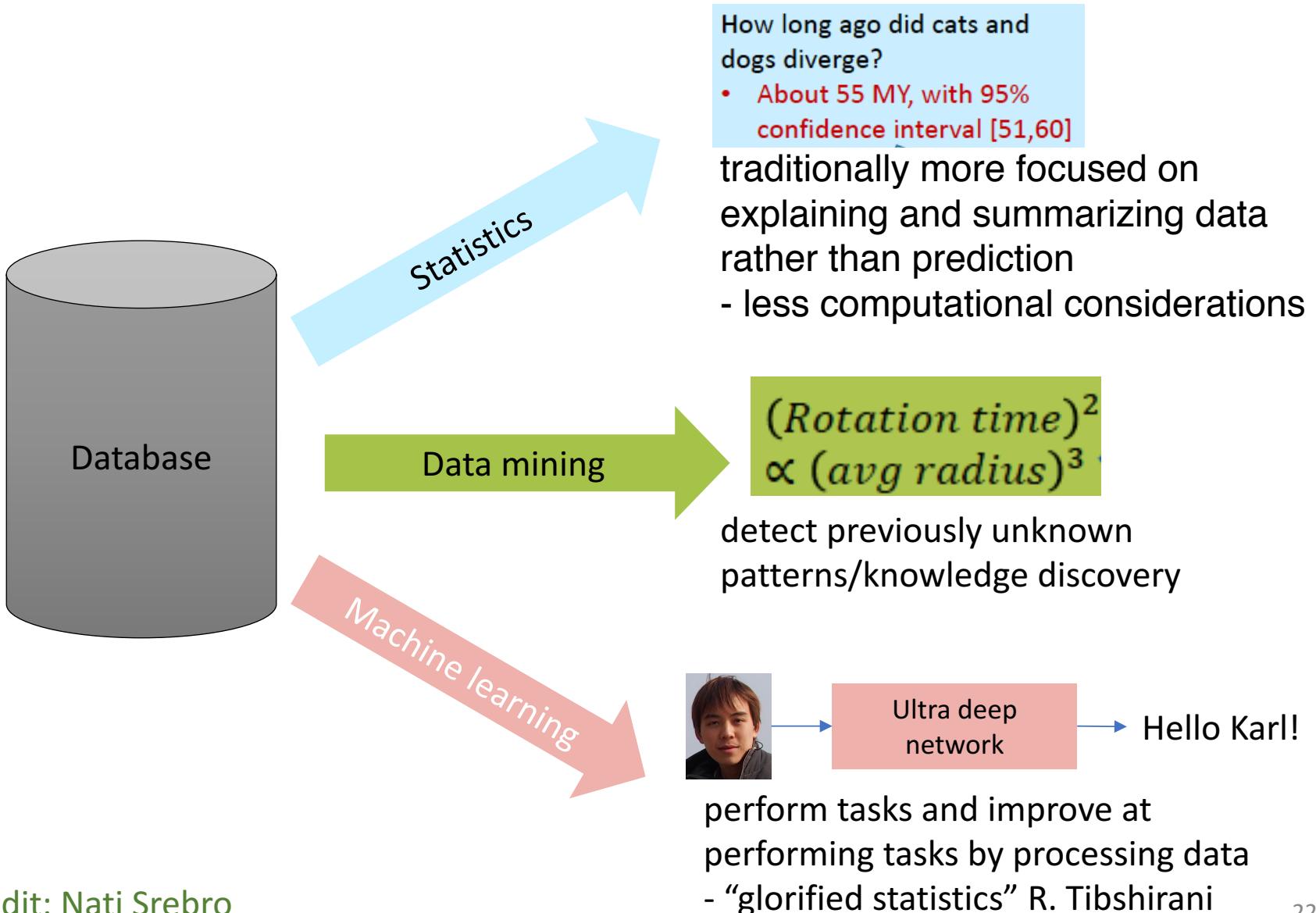
- Save human effort
 - automate boring tasks: surveillance, scanning license plates, driving!
 - automate dangerous tasks: robots for rescue operations, driving!
 - automate tasks requiring expert skills: translation
- Machines may perform tasks better than (expert) humans
 - detecting anomalies in radiology measurements
 - more accurate and consistent sensors!

Why learning?

- Machines can operate at scale and adapt
 - personalizing medicine, recommendations, etc.
 - learning based systems can adapt to new data
- Some systems are simply too complex to design manually
 - autonomous driving
 - weather prediction
- We have a lot of data and lot of computation
- Automate programming/automation



Data driven fields



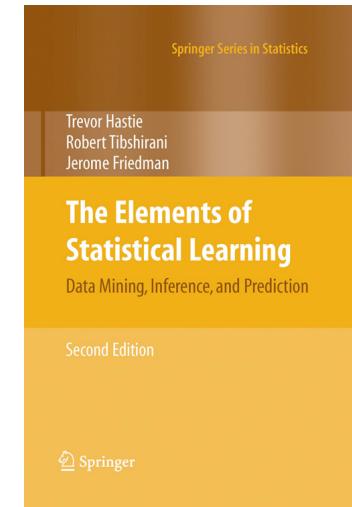
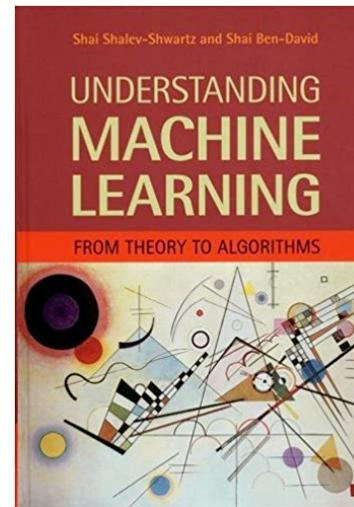
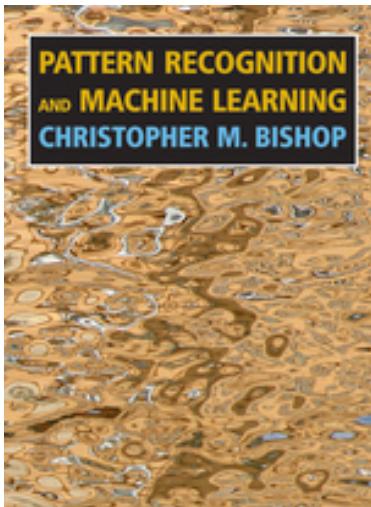
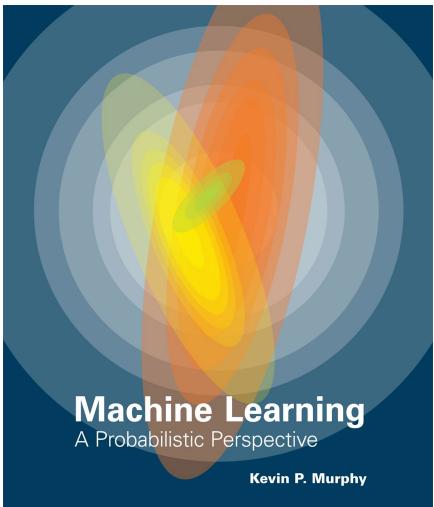
Goals of the course

- Understand basic models in machine learning, mainly supervised learning
- Familiarize with mathematical tools to formalize machine learning problems
- Gain hands-on experience in systematically implementing, diagnosing, and testing machine learning models
- (Somewhat) understand how/when/why machine learning works

Course overview

- Week 1
 - Basic models for supervised learning
 - Regression: linear regression
 - Classification: logistic regression, maximum margin classifiers, generative models for classification, structured classification
- Week 2
 - Introduction to neural networks, design and algorithmic choices in neural networks
 - Unsupervised learning

Resources



Supervised learning

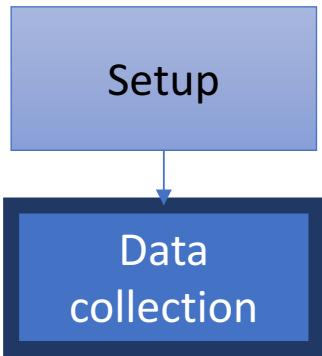
learn a function that maps input x to output y
from examples of (x, y)

Supervised learning: problem setup

Setup

- Input $x \in \mathcal{X}$
 - Example 1. meteorological measurements – \mathcal{X} numeric data in \mathbb{R}^d + text_reports
 - Example 2. RGB pixel values of an image $\mathcal{X} = \mathbb{R}^{h \times w \times 3}$
- Output/label/target $y \in \mathcal{Y}$
 - Example 1. average temperature for tomorrow $\mathcal{Y} = \mathbb{R}$
 - Example 2. face or no-face $\mathcal{Y} = \{\text{face, noface}\}$ or $\mathcal{Y} = \{0,1\}$
- Goal: learn a function to predict y from x
 $f: \mathcal{X} \rightarrow \mathcal{Y}$

Supervised learning: data collection



- Collect a set of labeled examples
$$S = \{(\mathbf{x}^{(1)}, \mathbf{y}^{(1)}), (\mathbf{x}^{(2)}, \mathbf{y}^{(2)}), \dots, (\mathbf{x}^{(N)}, \mathbf{y}^{(N)})\}$$
 - Example 1. weather prediction -> historic data
 - Example 2. face detection -> ask humans to label images with faces
- Desirable properties
 - want examples to be accurate – minimize noise

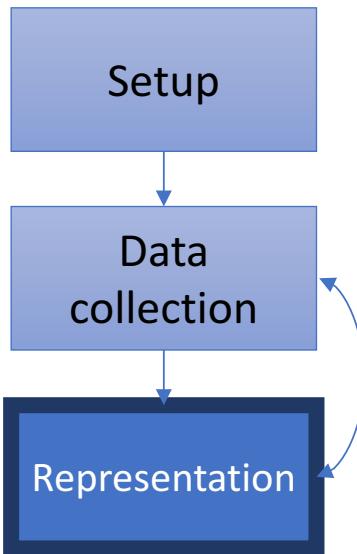
don't want  , noface

→ want examples to be diverse

don't want $\left\{ \left(\text{, face} \right), \left(\text{, face} \right), \left(\text{, face} \right) \right\}$

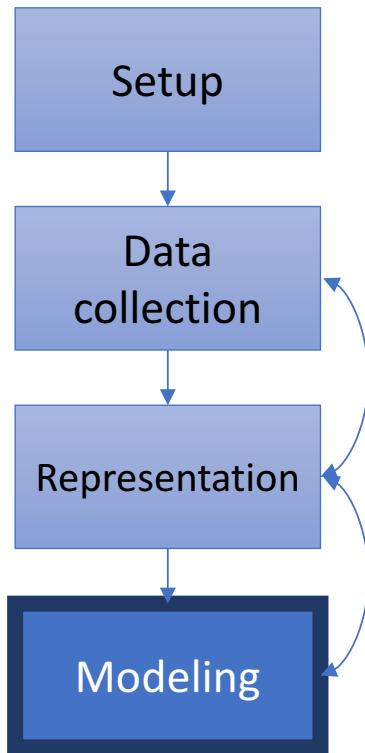
→ want many examples

Supervised learning: representation



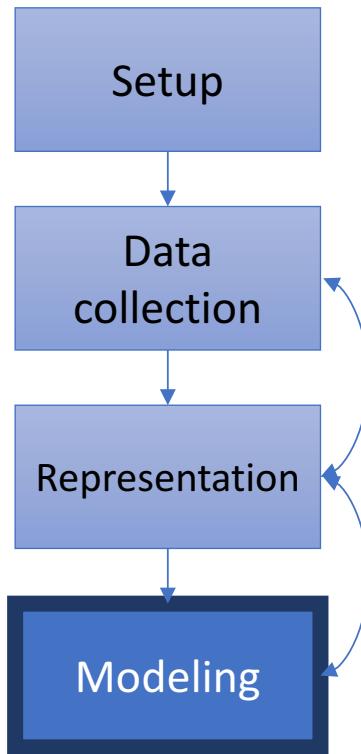
- Preprocess features $x \leftarrow \text{preprocess}(x)$
 - e.g., convert non-numeric data to numbers
 - what units to use for various measurements?
e.g., previous day temperature in $^{\circ}\text{C}$ or $^{\circ}\text{F}$?
does (should) it matter?
 - depends on next step!!
- Extract more sophisticated features
 $x \leftarrow \text{advanced_features}(x)$
 - e.g., pass raw pixels of an image through an edge detector
 - again depends on next step!!

Supervised learning: modeling



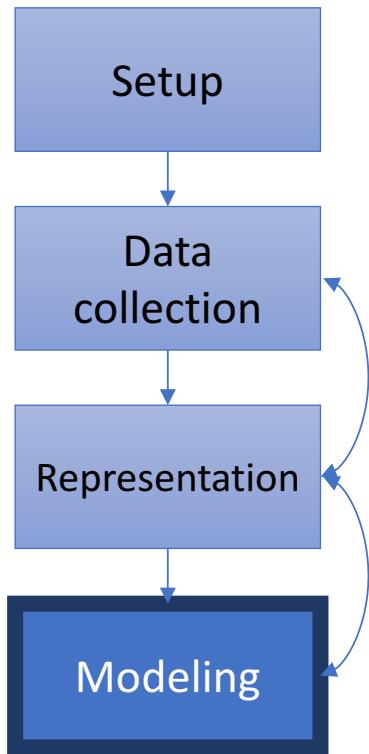
- Pick a set of candidates \mathcal{H} for $f: \mathcal{X} \rightarrow \mathcal{Y}$?
- Parameterize f by $w \in \mathbb{R}^P$, $\mathcal{H} = \{f_w: w \in \mathbb{R}^P\}$
then somehow pick “correct” w or “correct” f_w
(next step)
 - Example 1. $f(x)$ is a linear combination of numeric measurements in x
$$f_w(x) = w_1 x_1 + w_2 x_2 + \dots + w_d x_d = \mathbf{w} \cdot \mathbf{x}_{\text{num}}$$
 - Example 2. $f(x)$ is a linear combination of edge features
$$f_w(x) = \mathbf{w} \cdot \mathbf{x}_{\text{edge}} = \langle \mathbf{w}, \mathbf{x}_{\text{edge}} \rangle = \mathbf{w}^\top \mathbf{x}_{\text{edge}}$$

Supervised learning: modeling



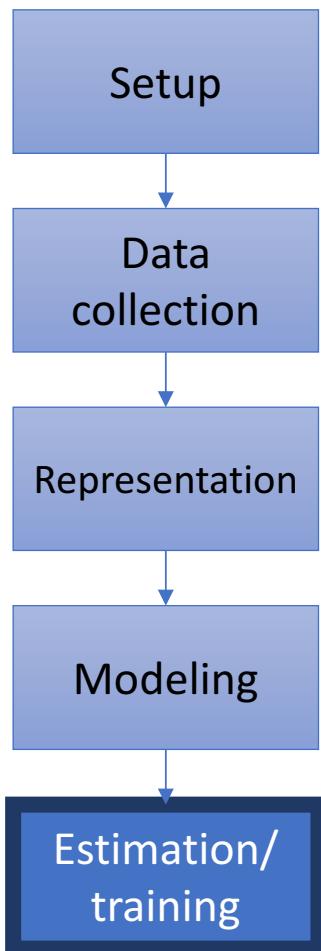
- Pick a set of candidates \mathcal{H} for $f: \mathcal{X} \rightarrow \mathcal{Y}$?
- Parameterize f by $w \in \mathbb{R}^P$ $\mathcal{H} = \{f_w: w \in \mathbb{R}^P\}$
then somehow pick “correct” w or “correct” f_w (next step)
 - Example 1. $f(x)$ is a linear combination of numeric measurements in x
$$f_w(x) = w_1 x_1 + w_2 x_2 + \dots w_d x_d = w \cdot x_{\text{num}}$$
 - Example 2. $f(x)$ is a linear combination of edge features
$$f_w(x) = w \cdot x_{\text{edge}} = \langle w, x_{\text{edge}} \rangle = w^\top x_{\text{edge}}$$
- Can also be non-parametric,
 - e.g., to compute $f(x)$ first find “closest” $x^{(i)}$ in training data and then return corresponding label $y^{(i)}$

Supervised learning: modeling



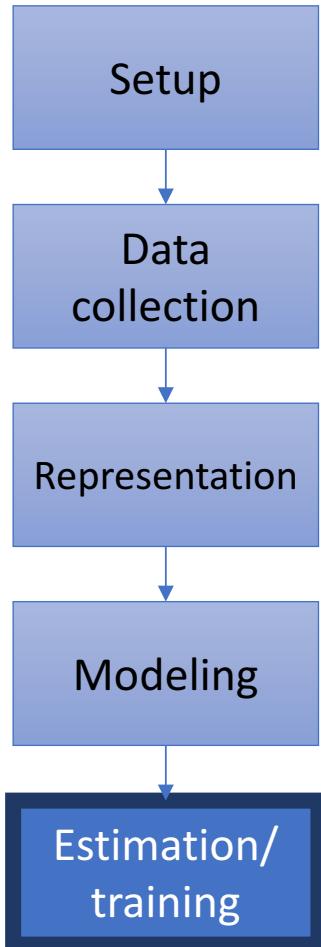
- Pick a set of candidates \mathcal{H} for $f: \mathcal{X} \rightarrow \mathcal{Y}$?
- Parameterize f by $w \in \mathbb{R}^P$ $\mathcal{H} = \{f_w: w \in \mathbb{R}^P\}$
then somehow pick “correct” w or “correct” f_w (next step)
 - Example 1. $f(x)$ is a linear combination of numeric measurements in x
$$f_w(x) = w_1 x_1 + w_2 x_2 + \dots w_d x_d = w \cdot x_{\text{num}}$$
 - Example 2. $f(x)$ is a linear combination of edge features
$$f_w(x) = w \cdot x_{\text{edge}} = \langle w, x_{\text{edge}} \rangle = w^T x_{\text{edge}}$$
- Can also be non-parametric,
 - e.g., to compute $f(x)$ first find “closest” $x^{(i)}$ in training data and then return corresponding label $y^{(i)}$
- Some main considerations
 - Does $\mathcal{H} = \{f_w: w \in \mathbb{R}^P\}$ capture what “correct” f is?
 - Do we have enough data to find the “correct” w ? (next step)
 - What is the cost of estimating “correct” w ? (next step)
 - Can $f(x) \notin \mathcal{Y}$? what if f returns -1 for no-faces instead of 0 ?

Supervised learning: training



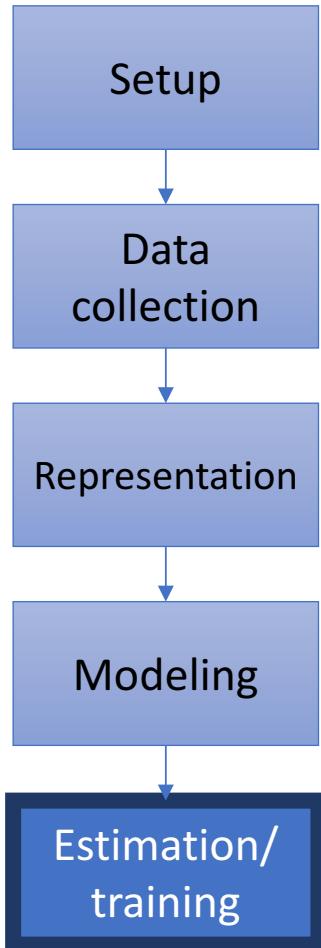
- Pick/estimate the “correct” \hat{f} from \mathcal{H}
- Option 1: Find $\hat{f} \in \mathcal{H}$ such that for all training examples $(x^{(i)}, y^{(i)}) \in S$ we get correct prediction,
$$\hat{f}(x^{(i)}) = y^{(i)} \text{ for } i = 1, 2, \dots, N$$
What happens if there is no such $f \in \mathcal{H}$? or many such $f \in \mathcal{H}$?

Supervised learning: training



- Pick/estimate the “correct” \hat{f} from \mathcal{H}
- Option 1: Find $\hat{f} \in \mathcal{H}$ such that for all training examples $(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}) \in S$ we get correct prediction,
$$\hat{f}(\mathbf{x}^{(i)}) = \mathbf{y}^{(i)} \text{ for } i = 1, 2, \dots, N$$
What happens if there is no such $f \in \mathcal{H}$? or many such $f \in \mathcal{H}$?
- Option 2: Find $\hat{f} \in \mathcal{H}$ that “fits/explains” most of data?
 - find \hat{f} that minimizes some “loss” ℓ between prediction $f(\mathbf{x}^{(i)})$ and $\mathbf{y}^{(i)}$
$$\hat{f} = \operatorname{argmin}_{f \in \mathcal{H}} \sum_{i=1}^N \ell(f(\mathbf{x}^{(i)}), \mathbf{y}^{(i)}) \text{ or } \hat{w} = \operatorname{argmin}_{w \in \mathbb{R}^d} \sum_{i=1}^N \ell(f_w(\mathbf{x}^{(i)}), \mathbf{y}^{(i)})$$

Supervised learning: training



- Pick/estimate the “correct” \hat{f} from \mathcal{H}
- Option 1: Find $\hat{f} \in \mathcal{H}$ such that for all training examples $(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}) \in S$ we get correct prediction,

$$\hat{f}(\mathbf{x}^{(i)}) = \mathbf{y}^{(i)} \text{ for } i = 1, 2, \dots, N$$

What happens if there is no such $f \in \mathcal{H}$? or many such $f \in \mathcal{H}$?

- Option 2: Find $\hat{f} \in \mathcal{H}$ that “fits/explains” most of data?
 - find \hat{f} that minimizes some “loss” ℓ between prediction $f(x_i)$ and y_i

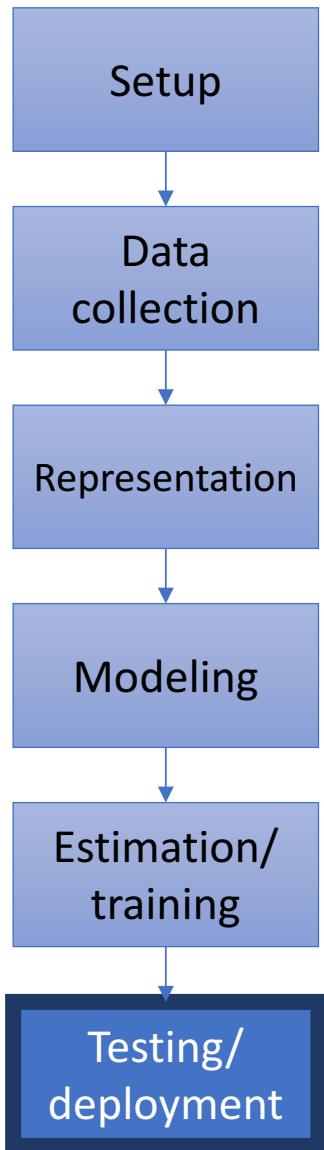
$$\hat{f} = \operatorname{argmin}_{f \in \mathcal{H}} \sum_{i=1}^N \ell(f(\mathbf{x}^{(i)}), \mathbf{y}^{(i)}) \text{ or } \hat{w} = \operatorname{argmin}_{w \in \mathbb{R}^d} \sum_{i=1}^N \ell(f_w(\mathbf{x}^{(i)}), \mathbf{y}^{(i)})$$

- what loss ℓ to use?
- how to solve the minimization? – computational considerations
- Example 1: $\ell(f(\mathbf{x}^{(i)}), \mathbf{y}^{(i)}) = |f(\mathbf{x}^{(i)}) - \mathbf{y}^{(i)}|$ or $(f(\mathbf{x}^{(i)}) - \mathbf{y}^{(i)})^2$

$$\text{◦ Example 2: } \ell(f(\mathbf{x}^{(i)}), \mathbf{y}^{(i)}) = \begin{cases} 0 & \text{if } f(\mathbf{x}^{(i)}) = \mathbf{y}^{(i)} \\ 1 & \text{if } f(\mathbf{x}^{(i)}) \neq \mathbf{y}^{(i)} \end{cases}$$

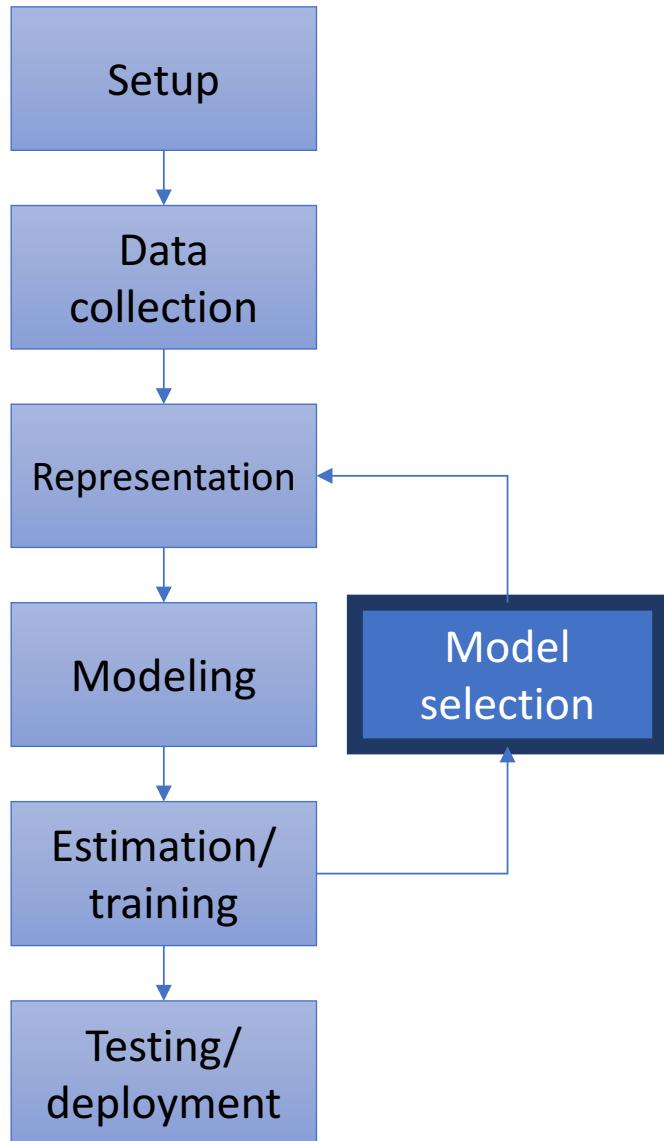
→ Empirical Risk Minimization (ERM)

Supervised learning: testing



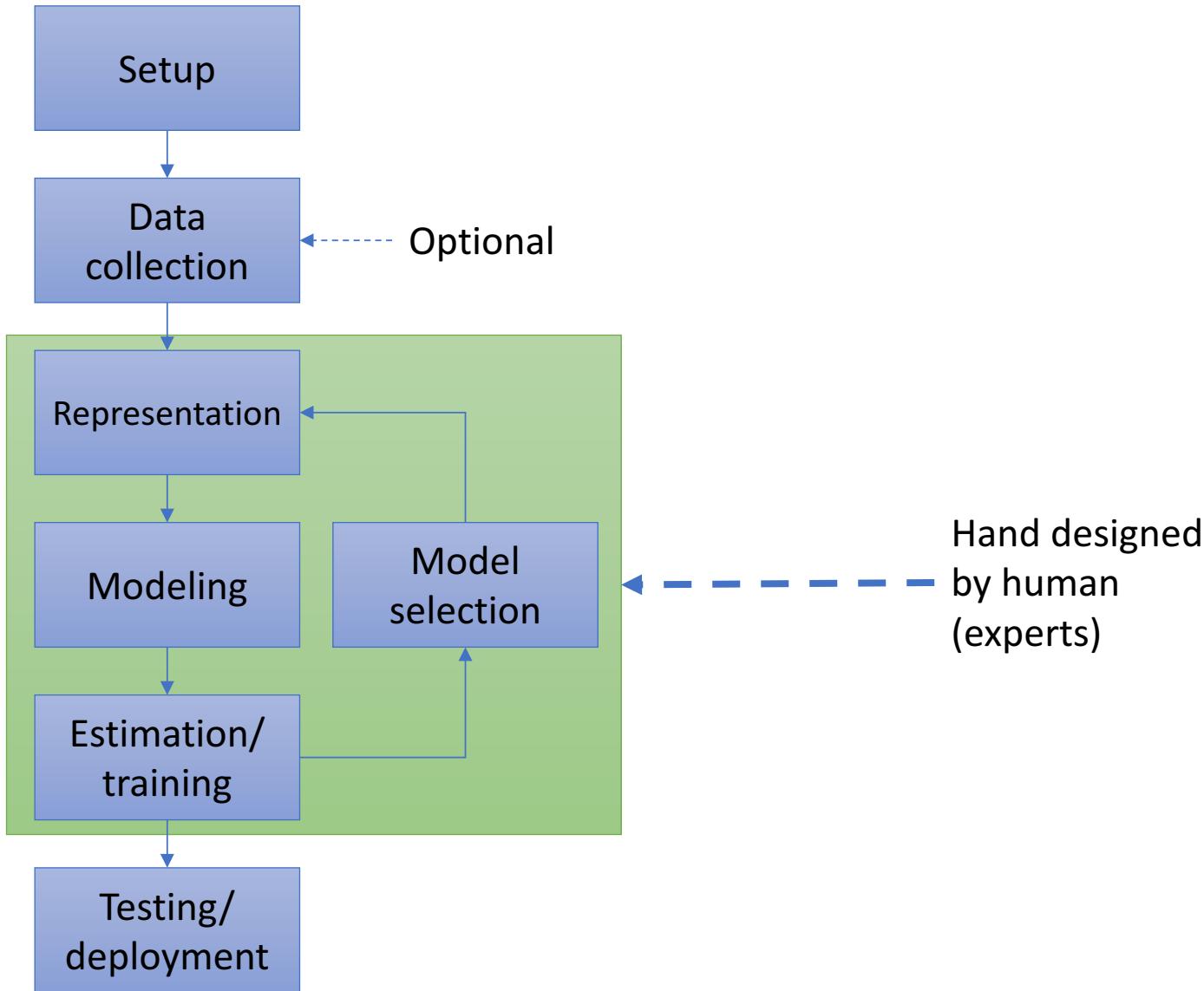
- Test how well \hat{f} performs the task
- Can we check $\sum_i \ell(\hat{f}(\mathbf{x}^{(i)}), y^{(i)})$?
 - e.g., $\left\{ \left(\begin{array}{c} \text{[image of Karl]} \\ \text{[image of Karl]} \end{array}, \text{face} \right), \left(\begin{array}{c} \text{[image of Karl]} \\ \text{[image of Karl]} \end{array}, \text{face} \right), \left(\begin{array}{c} \text{[black box]} \\ \text{[black box]} \end{array}, \text{noface} \right) \right\}$
and \hat{f} always outputs “noface” for any non-Karl image
- What we want is how well \hat{f} performs on new examples!
- Deploy in real world and see how well \hat{f} performs
- Collect more data S_{test} but do not use for training
- Split original data and do training on only S_{train}

Supervised learning: Model selection

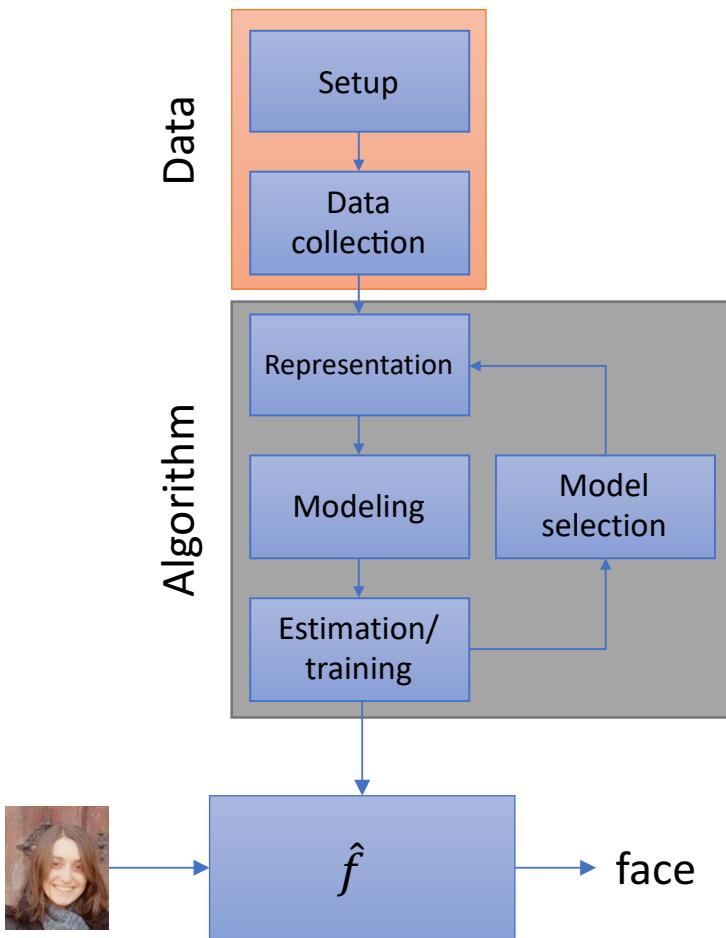


- What if \hat{f} performs poorly on S_{test} ?
 - Can we again use same S_{test} ?
- Can we “detect” how well \hat{f} will perform while testing?
- Mimic “testing” during model selection
 - Keep a small number of training example to mimic the testing process – validation or development data
 - Split S into S_{train} , S_{val} and S_{test}
 - Train on S_{train} , get mock-test score by checking on S_{val}
 - Change models and repeat
 - DO NOT use S_{test}

Rule based learning



Supervised learning – key questions



- **Data:** what kind of data can we get? how much data can we get?
- **Model:** what is the correct model for my data? – want to minimize the effort put into this question!
- **Training:** what resources - computation/memory - does the algorithm need to estimate the model \hat{f} ?
- **Testing:** how well will \hat{f} perform when deployed? what is the computational/memory requirement during deployment?

Other learning paradigms

- **Unsupervised learning (will cover some)**
 - Only examples of x are seen but no y
 - There might not be any specific y
 - Goal is to make sense of the data
- **Hybrid learning**
 - semi-supervised: small number of (x, y) examples are available, but lot more examples of just x
 - weakly supervised: (x, \tilde{y}) examples are available where \tilde{y} is a very noisy estimate of y , or has partial information, e.g., fruit instead of apple
- **Reinforcement learning**
 - Learning to plan and execute long term actions
 - e.g., navigate a maze, win a soccer match