

# Overview

CS 446 / ECE 449

2022-01-17 04:49:48 -0600 (d748537)

# Plan for today

- ▶ Course webpage: staff, policies, schedule.
- ▶ ML examples.
- ▶ ML math/coding background.
- ▶ ML setting and meta-algorithms.
- ▶ First ML technique: linear regression.

# What is machine learning?

Improving algorithms by fitting them to data.

## Application 1: image classification

- ▶ Birdwatcher takes photos of birds, organized by species.
- ▶ **Goal:** automatically recognize bird species in new photos.

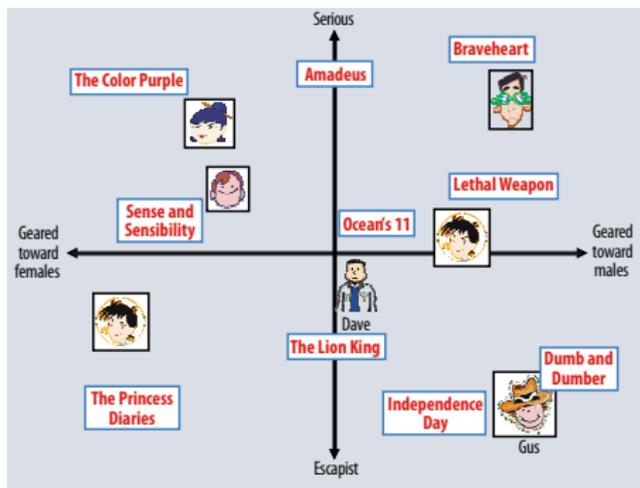


Indigo bunting

- ▶ **Why ML:** variation in lighting, occlusions, morphology.

## Application 2: recommender system

- ▶ Netflix users watch movies and provide ratings.
- ▶ **Goal:** predict user's rating of unwatched movie.
- ▶ (**Real goal:** keep users paying customers.)
- ▶ (**Real effect:** reinforce stereotypes found in the data?)

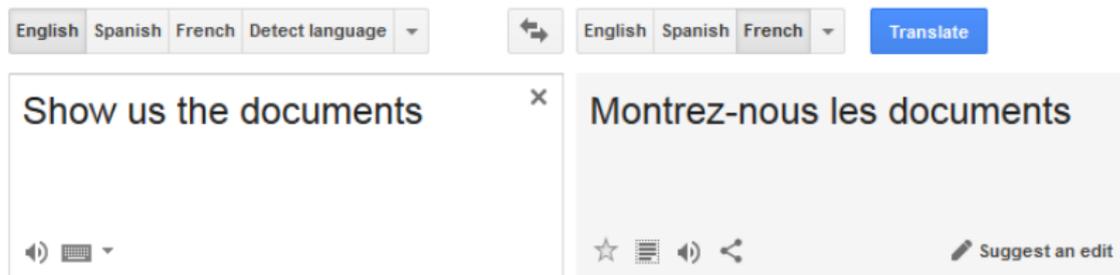


(Image credit: Koren, Bell, and Volinsky, 2009.)

- ▶ **Why ML:** easily adapt to and leverage movie attributes, viewer attributes, viewer relationships, etc.

## Application 3: machine translation

- ▶ Linguists provide translations of all English language books into French, sentence-by-sentence.
- ▶ **Goal:** translate any English sentence into French.



**Note:** the text-to-speech is via ML (~~recurrent network~~ transformer).

- ▶ **Why ML?** Not only avoid hard-coding many rules, but also capture idiom and other nuances.

# Application 4: chess

- ▶ Chess enthusiasts construct a large corpus of chess games.  
(Or: start with nothing, play random games!)
- ▶ **Goal:** Win chess games.

The screenshot shows a chess game on lichess.org. The board is in a checkmate position. White's king is on g1, and Black's king is on g8. White's queen is on g6, and Black's queen is on e4. The game is over with a score of 1-0, indicating a checkmate. The interface includes a sidebar with game details, a main board area, and a move list on the right. The move list shows the following moves: 54... ♖b7 55. ♖c4 ♜b8 [...], 55. ♖c4 +13.3 f5?! #21, 56. ♜g1 +5.4 ♜f2+?! #10, 56... ♜b7 57. ♜g8, 57. ♜e2 #9 d5 #8, 58. ♜a6+ #7 ♜b8 #7, 59. ♜g8+ #6 ♜c7 #6, 60. ♜c8+ #5 ♜d7 #5, 61. ♜c6+ #4 ♜e7 #4, 62. ♜e8+ #3 ♜f7 #3, 63. ♜e6+ #2 ♜g7 #2, 64. ♜g8+ #1 ♜h7 #1, 65. ♜g6#.

Game details: 5+2 • Rated • Blitz 3 weeks ago, 36 stars. Opponents: BOT Stockfish10Chess (25837) +21, BOT LeelaChess (27607) -43. Result: Checkmate • White is victorious.

Game statistics: 2 Spectators GM penguimgim1, Anonymous. Stockfish10Chess +214, 0 Inaccuracies, 0 Mistakes, 0 Blunders, 11 Average centipawn loss. A button to 'Learn from your mistakes' is visible.

- ▶ **Why ML?** Avoid hard-coding evaluation; magically interpolate between observed positions, as humans do.



## Math.

- ▶ Linear algebra (e.g., null spaces; eigendecomposition; SVD...).
- ▶ Basic probability and statistics (e.g., variance of a random variable).
- ▶ Multivariable calculus (e.g., gradient of  $\|\mathbf{A}\mathbf{w} - \mathbf{b}\|^2$  wrt  $\mathbf{w}$ ).
- ▶ Basic proof writing (e.g., prove  $\mathbf{A}^T\mathbf{A}$  is positive semi-definite).

## Math.

- ▶ Linear algebra (e.g., null spaces; eigendecomposition; SVD...).
- ▶ Basic probability and statistics (e.g., variance of a random variable).
- ▶ Multivariable calculus (e.g., gradient of  $\|A\mathbf{w} - \mathbf{b}\|^2$  wrt  $\mathbf{w}$ ).
- ▶ Basic proof writing (e.g., prove  $A^T A$  is positive semi-definite).

## Coding.

- ▶ `python3`. It's slow, but often computation will be inside fast libraries.
- ▶ `numpy`, an easy-to-use numeric library.
- ▶ `pytorch`, a numeric library with gpu support, auto-differentiation, and deep learning helpers.

**My opinion.** `pytorch` is one of the nicest libraries I've ever used, for anything. I use it for much more than deep learning.

```
>>> import numpy
>>> import torch

>>> 3 / 2
1.5
>>> 3 // 2
1

>>> A = torch.randn(5,5)
>>> b = torch.randn(5,1)
>>> x = torch.gels(b, A)[0]
>>> (A @ x - b).norm()
tensor(3.7985e-06)
```

**See also:** pytorch tutorial (lecture 9).

## pytorch on gpu is easy

```
>>> import torch

>>> A = torch.randn(5, 5)
>>> b = torch.randn(5)
>>> (A @ b).norm()
tensor(4.7746)

>>> device = torch.device("cuda:0")
>>> (A.to(device) @ b.to(device)).norm()
tensor(4.7746, device='cuda:0')

>>> (A.to(device) @ b).norm()
Traceback (most recent call last):
  File "<stdin>", line 1, in <module>
RuntimeError: Expected object of type
torch.cuda.FloatTensor but found type
torch.FloatTensor for argument #2 'vec'
```

**Note.** Homeworks will be graded in a gpu-free container!

# Simplest setting: supervised learning

## Supervised learning.

1. Receive **training set**  $((\mathbf{x}_i, y_i))_{i=1}^n$ , where each  $\mathbf{x}_i$  is an **input** or **covariate**, and each  $y_i$  is a **label** or **target**.  
("Supervision": something gives us the labels!)
2. Algorithmically choose a predictor  $f$  via the **training set** so that  $f(\mathbf{x}) \approx y$  on future examples.

# Simplest setting: supervised learning

## Supervised learning.

1. Receive **training set**  $((\mathbf{x}_i, y_i))_{i=1}^n$ , where each  $\mathbf{x}_i$  is an **input** or **covariate**, and each  $y_i$  is a **label** or **target**.  
("Supervision": something gives us the labels!)
2. Algorithmically choose a predictor  $f$  via the **training set** so that  $f(\mathbf{x}) \approx y$  on future examples.

## Other settings.

- ▶ **Unsupervised learning**: find structure in  $(\mathbf{x}_i)_{i=1}^n$  (no labels!).
- ▶ **Time series** label of  $\mathbf{x}_i$  depends on  $(\mathbf{x}_{i-1}, \mathbf{x}_{i-2}, \dots)$ .
- ▶ **Reinforcement learning**: predictions/outputs affect future state (e.g., driving a car).

# What are the difficulties?

Consider supervised learning (simplest setting):  
learn  $f : X \rightarrow Y$  from  $((\mathbf{x}_i, y_i))_{i=1}^n$ .

# What are the difficulties?

Consider supervised learning (simplest setting):  
learn  $f : X \rightarrow Y$  from  $((\mathbf{x}_i, y_i))_{i=1}^n$ .

- ▶ How to **encode data** for the algorithm?

# What are the difficulties?

Consider supervised learning (simplest setting):

learn  $f : X \rightarrow Y$  from  $((\mathbf{x}_i, y_i))_{i=1}^n$ .

- ▶ How to **encode data** for the algorithm?
- ▶ How to **clean/improve/augment** data?

# What are the difficulties?

Consider supervised learning (simplest setting):

learn  $f : X \rightarrow Y$  from  $((\mathbf{x}_i, y_i))_{i=1}^n$ .

- ▶ How to **encode data** for the algorithm?
- ▶ How to **clean/improve/augment** data?
- ▶ How to choose **structure/model** for  $f$ ?

# What are the difficulties?

Consider supervised learning (simplest setting):  
learn  $f : X \rightarrow Y$  from  $((\mathbf{x}_i, y_i))_{i=1}^n$ .

- ▶ How to **encode data** for the algorithm?
- ▶ How to **clean/improve/augment** data?
- ▶ How to choose **structure/model** for  $f$ ?
- ▶ How to algorithmically **fit**  $f$  to data?

# What are the difficulties?

Consider supervised learning (simplest setting):

learn  $f : X \rightarrow Y$  from  $((\mathbf{x}_i, y_i))_{i=1}^n$ .

- ▶ How to **encode data** for the algorithm?
- ▶ How to **clean/improve/augment** data?
- ▶ How to choose **structure/model** for  $f$ ?
- ▶ How to algorithmically **fit**  $f$  to data?
- ▶ How to ensure  $f$  does not **overfit**, meaning it is good on future predictions, and not just on  $((\mathbf{x}_i, y_i))_{i=1}^n$ ?

# “pytorch meta-algorithm”

1. **Clean/augment data** (lecture 10?).
2. **Pick model/architecture** (many lectures).
3. **Pick a loss function** measuring model fit to data (lectures 2-4, 6).
4. **Run a gradient descent variant** to fit model to data (many lectures).
5. Tweak 1-4 until **training error** is small.
6. Tweak 1-5, **possibly reducing model complexity**, until **testing error** is small (lectures 4, 6, 13).

# “pytorch meta-algorithm”

1. **Clean/augment data** (lecture 10?).
2. **Pick model/architecture** (many lectures).
3. **Pick a loss function** measuring model fit to data (lectures 2-4, 6).
4. **Run a gradient descent variant** to fit model to data (many lectures).
5. Tweak 1-4 until **training error** is small.
6. Tweak 1-5, **possibly reducing model complexity**, until **testing error** is small (lectures 4, 6, 13).

**Is that all of ML?**

# “pytorch meta-algorithm”

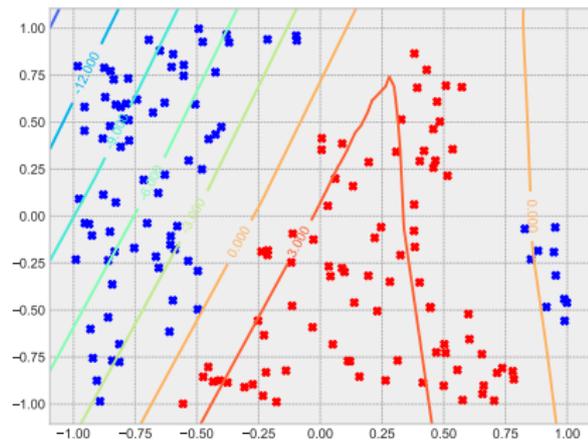
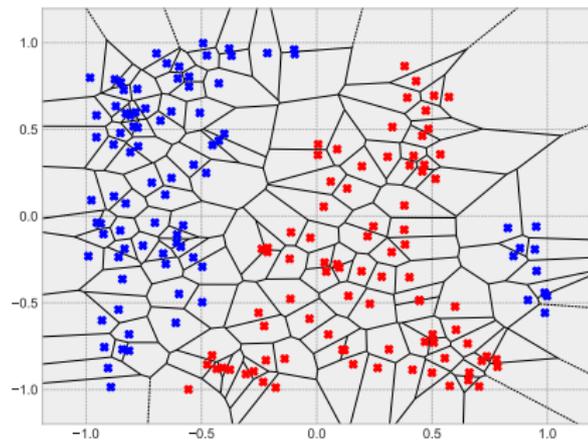
1. **Clean/augment data** (lecture 10?).
2. **Pick model/architecture** (many lectures).
3. **Pick a loss function** measuring model fit to data (lectures 2-4, 6).
4. **Run a gradient descent variant** to fit model to data (many lectures).
5. Tweak 1-4 until **training error** is small.
6. Tweak 1-5, **possibly reducing model complexity**, until **testing error** is small (lectures 4, 6, 13).

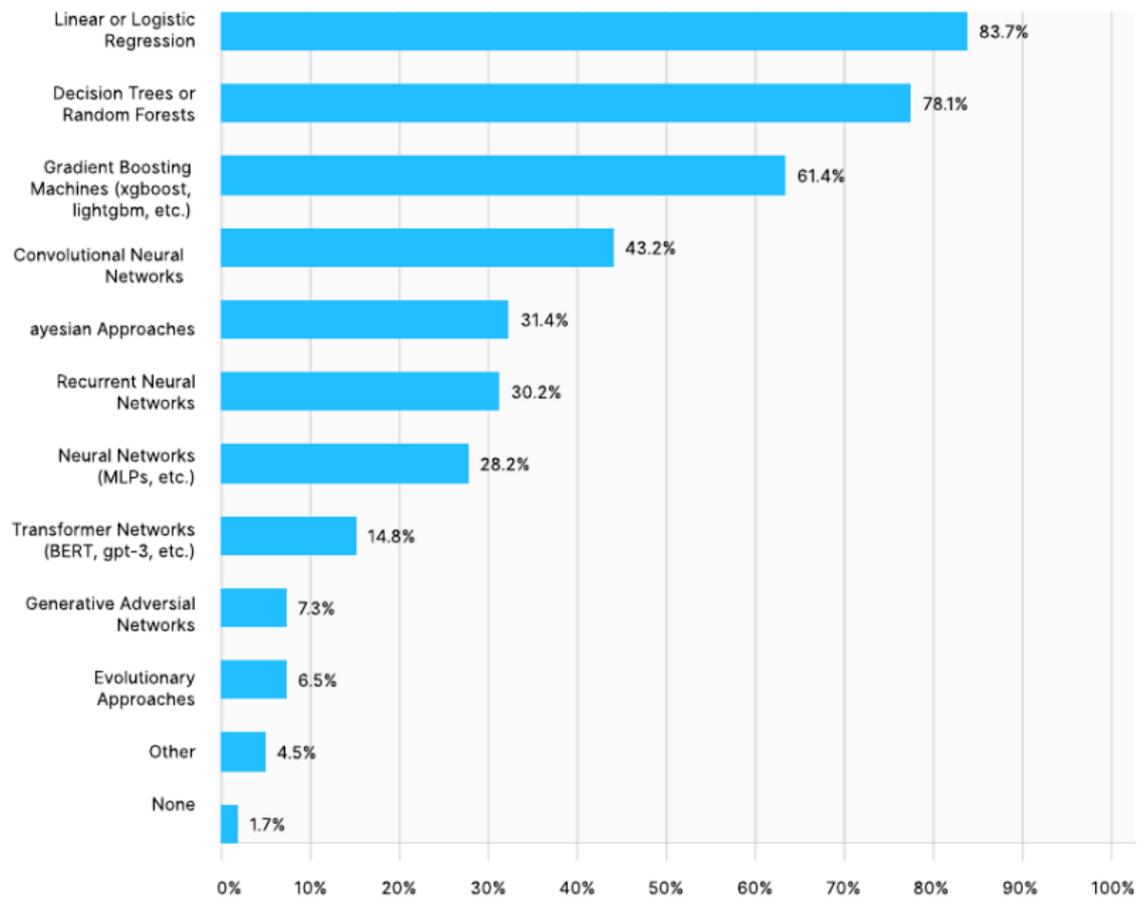
**Is that all of ML?**

**No**, but these days it's much of it!

# What is “model/architecture choice”?

Should we use 1-nearest-neighbor or a 2-layer ReLU network?





(From kaggle 2020 survey.)

# Linear regression — basic setup

1. Start from **training data**  $((\mathbf{x}_i, y_i))_{i=1}^n$ , with  $\mathbf{x}_i \in \mathbb{R}^d$  and  $y_i \in \mathbb{R}$ .
2. **Model** is a **linear predictor**: pick  $\mathbf{w} \in \mathbb{R}^d$  with

$$\mathbf{x}_i \mapsto \mathbf{w}^\top \mathbf{x}_i =: \hat{y}_i \approx y_i.$$

3. **Loss function**  $\ell$  is **squared loss**  $\ell_{\text{sq}}$  (standard regression loss):

$$\ell_{\text{sq}}(\mathbf{w}^\top \mathbf{x}_i, y_i) = \frac{1}{2}(\mathbf{w}^\top \mathbf{x}_i - y_i)^2.$$

We will minimize the **empirical risk** (average loss over training examples):

$$\hat{\mathcal{R}}(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \ell_{\text{sq}}(\mathbf{w}^\top \mathbf{x}_i, y_i) = \frac{1}{2n} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}_i - y_i)^2.$$

Convenient form using matrices:

$$\hat{\mathcal{R}}(\mathbf{w}) = \frac{1}{2n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 \quad \text{where } \mathbf{X} := \begin{bmatrix} \leftarrow & \mathbf{x}_1^\top & \rightarrow \\ & \vdots & \\ \leftarrow & \mathbf{x}_n^\top & \rightarrow \end{bmatrix} \in \mathbb{R}^{n \times d}.$$

### 3. Minimize the empirical risk

$$\widehat{\mathcal{R}}(\mathbf{w}) = \frac{1}{2n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 \quad \text{where } \mathbf{X} := \begin{bmatrix} \leftarrow & \mathbf{x}_1^\top & \rightarrow \\ & \vdots & \\ \leftarrow & \mathbf{x}_n^\top & \rightarrow \end{bmatrix} \in \mathbb{R}^{n \times d}.$$

### 4. Basic method: gradient descent. Set $\mathbf{w}_0 = 0$ , and thereafter

$$\mathbf{w}_{i+1} := \mathbf{w}_i - \eta \nabla \widehat{\mathcal{R}}(\mathbf{w}_i) = \mathbf{w}_i - \frac{\eta}{n} \mathbf{X}^\top (\mathbf{X}\mathbf{w}_i - \mathbf{y}),$$

where  $\eta$  is a learning rate (step size).

```
w = torch.zeros(d)
for i in range(niters):
    w -= eta * X.T @ (X @ w - y) / n
```

# Summary for today

- ▶ Course webpage: staff, policies, schedule.
- ▶ ML examples.
- ▶ ML math/coding background.
- ▶ ML setting and meta-algorithms.
- ▶ First ML technique: linear regression.

(Appendix.)

- ▶ Murphy: chapter 1.