

Machine Learning I

60629A

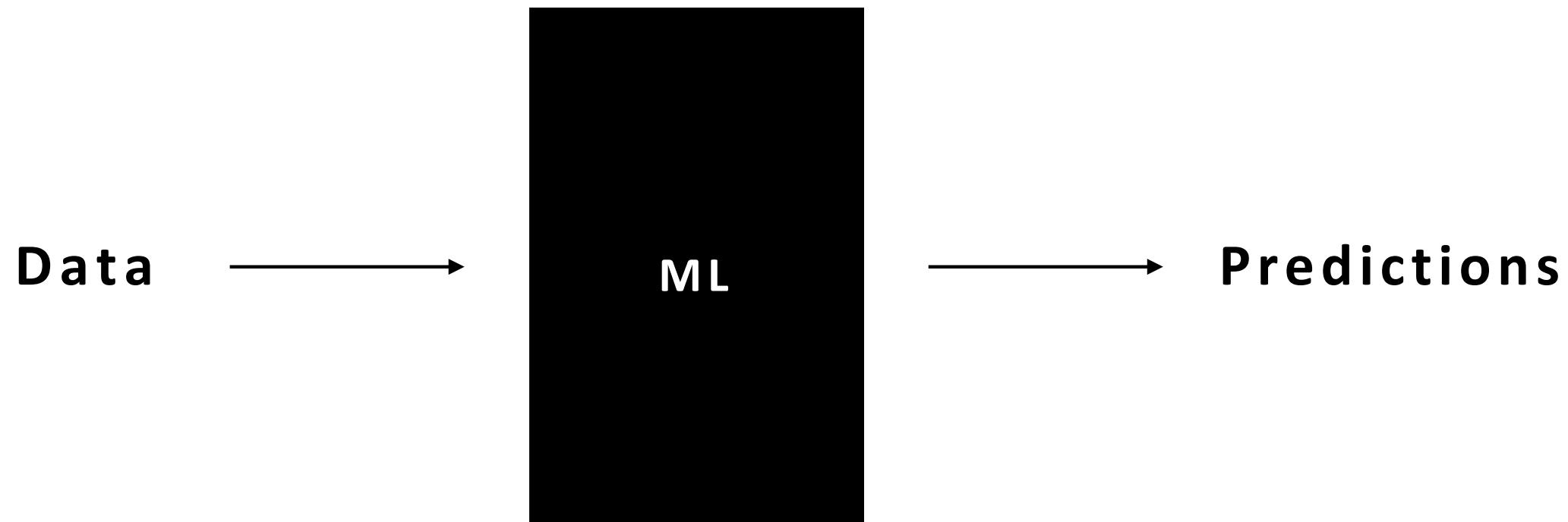
Week #1

Today

- **Introduction to machine learning**
- **The course (syllabus)**
- **Math review (probability + linear algebra)**

Machine Learning (ML)

- Science that studies statistical and computational aspects of modeling data for predictive purposes
- (Mostly) Empirical science



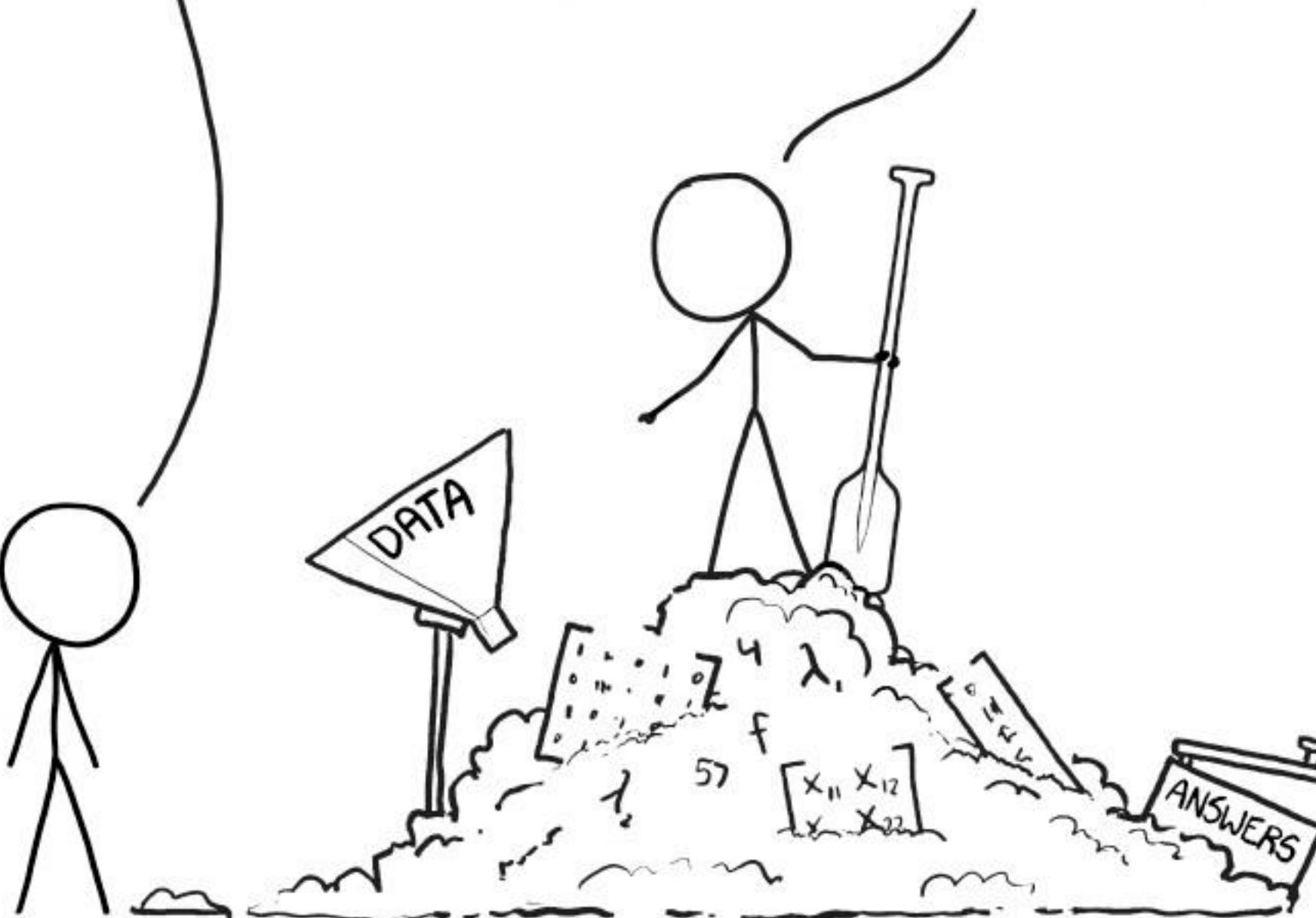
- **Task: Predict whether an image contains a malignant tumor**
- **Task: Predict the next movie a person should watch**

THIS IS YOUR MACHINE LEARNING SYSTEM?

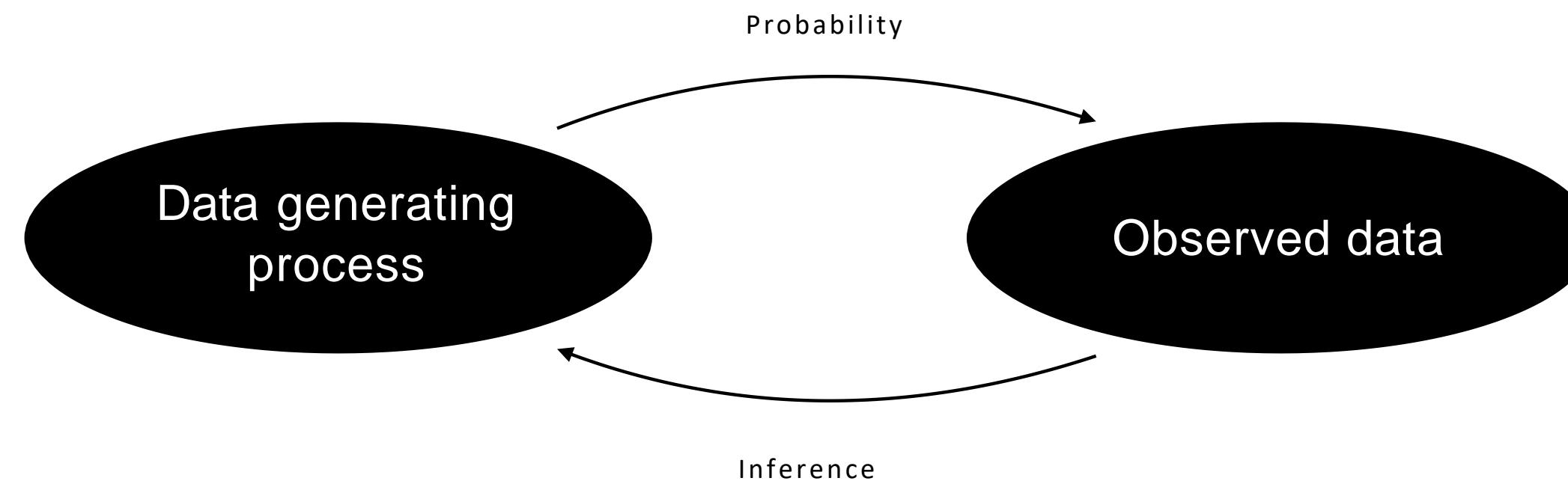
YUP! YOU POUR THE DATA INTO THIS BIG
PILE OF LINEAR ALGEBRA, THEN COLLECT
THE ANSWERS ON THE OTHER SIDE.

WHAT IF THE ANSWERS ARE WRONG?

JUST STIR THE PILE UNTIL
THEY START LOOKING RIGHT.



“Data analysis, machine learning and data mining are various names given to the practice of statistical inference, depending on the context.”



-Larry Wasserman in “All of Statistics: A Concise Course in Statistical Inference.”

What is the goal of ML?

- A bit of historical context
- When I started my PhD very few in ML talked about AI
- Recent ML makes progress toward “AI tasks”
 - Examples of AI tasks: translation, object recognition
 - In that context: create a machine with human-like capacities? Or a machine that can help humans?

- My (imperfect) view:

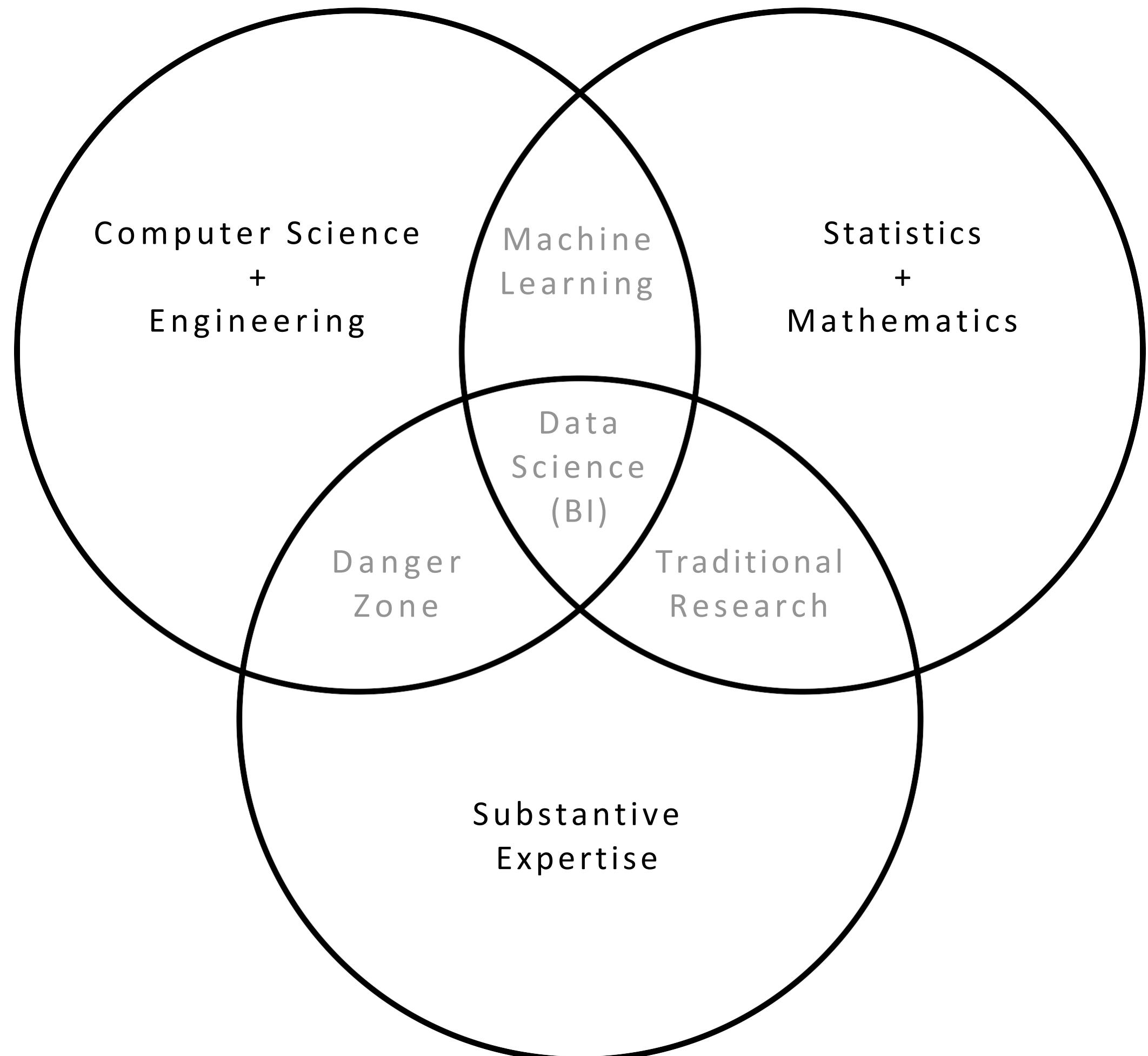
Understand data through predictive models

Understand the world through predictive models

How does ML relate
to other fields

Historical View

- **(Modern) Statistics:** ~1900
- **Machine Learning and Data Mining:** ~1960
- **Data Science:** ~2000



Attitudes in Machine Learning and Data Mining Versus Attitudes in Traditional Statistics

Despite these differences, there's a big overlap in problems addressed by machine learning and data mining and by traditional statistics. But attitudes differ...

Machine learning

No settled philosophy or widely accepted theoretical framework.

Willing to use *ad hoc* methods if they seem to work well (though appearances may be misleading).

Emphasis on automatic methods with little or no human intervention.

Methods suitable for many problems.

Heavy use of computing.

Traditional statistics

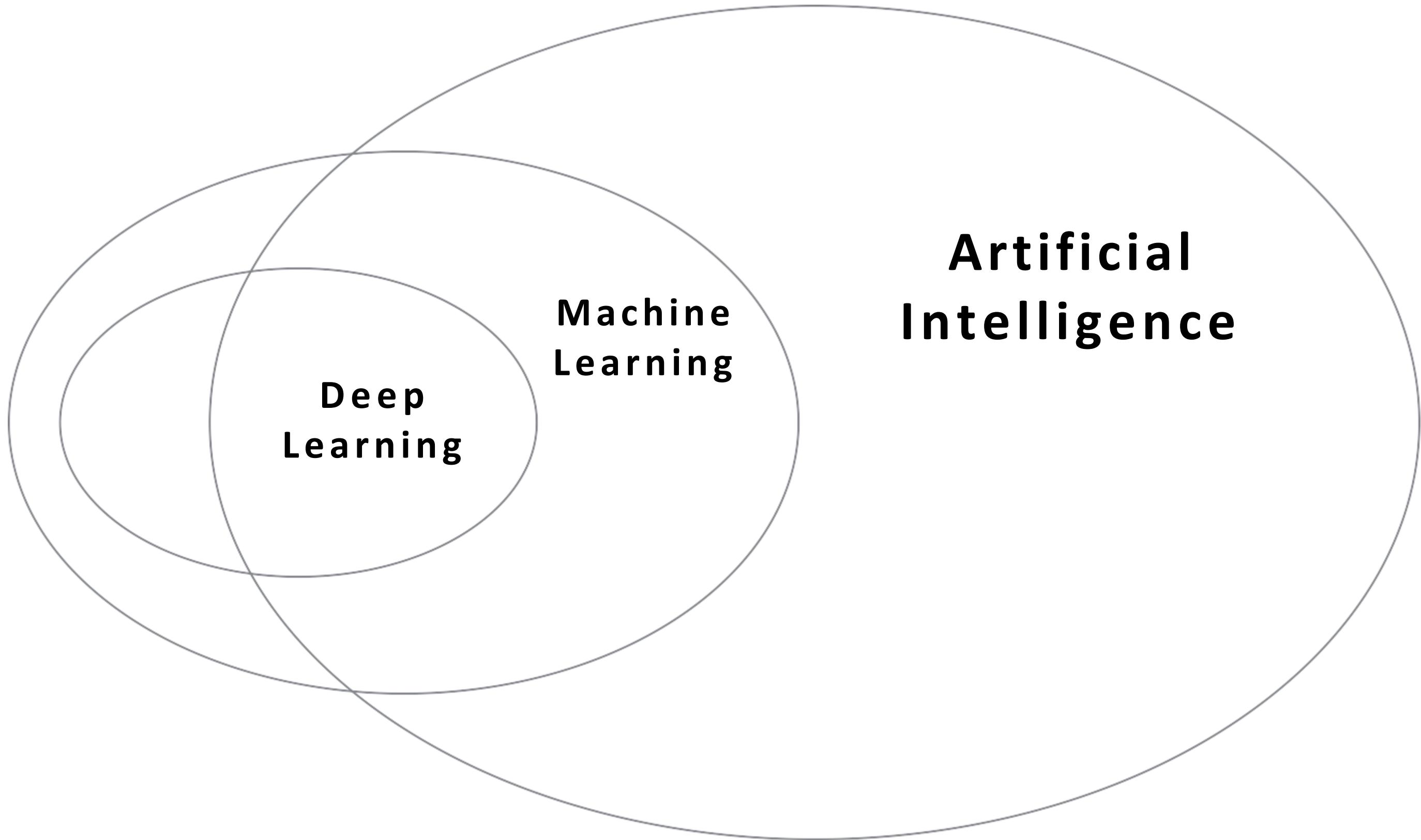
Classical (frequentist) and Bayesian philosophies compete.

Reluctant to use methods without some theoretical justification (even if the justification is actually meaningless).

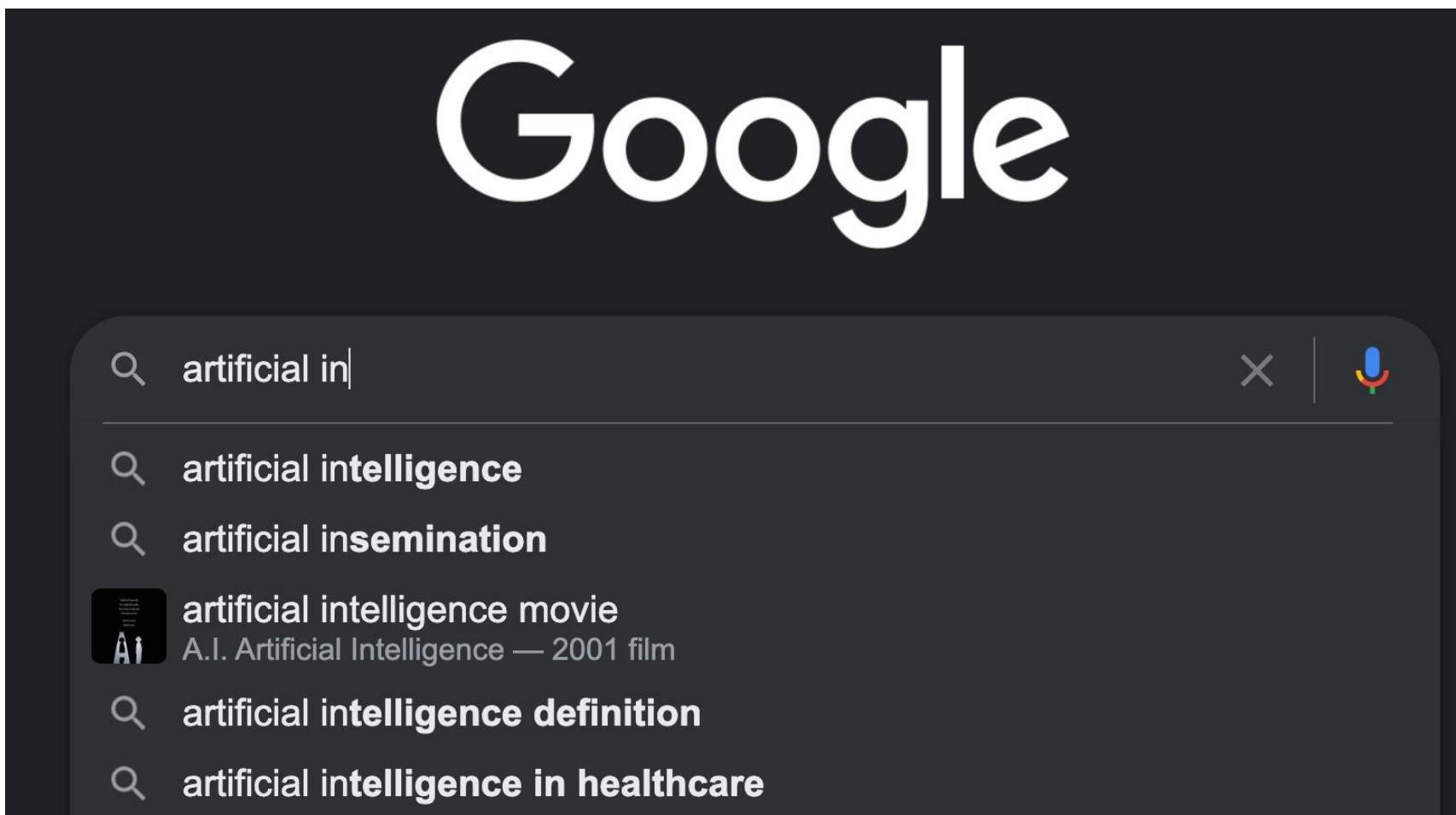
Emphasis on use of human judgement assisted by plots and diagnostics.

Models based on scientific knowledge.

Originally designed for hand-calculation, but computing is now very important.



Applications of ML



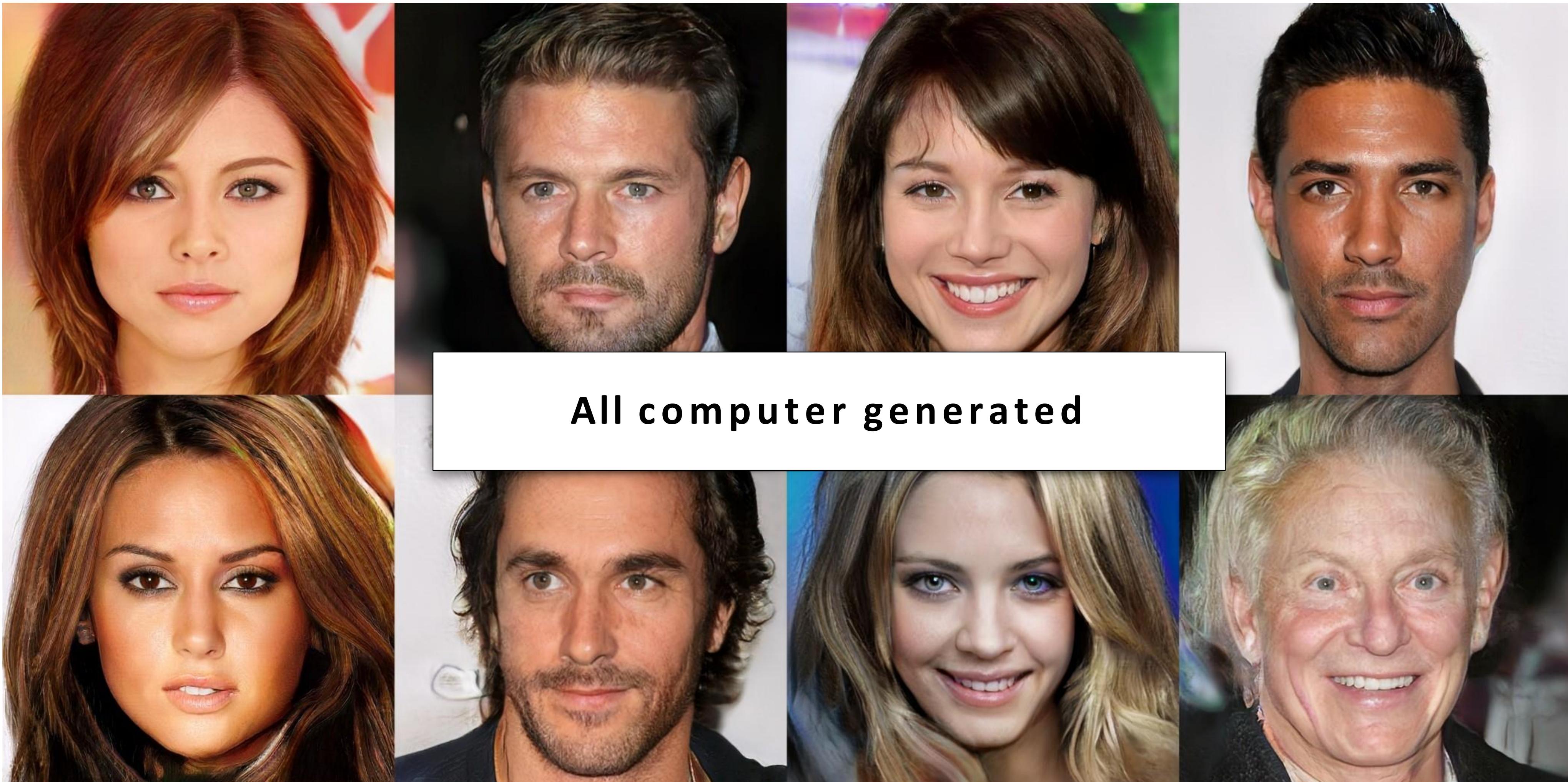
A screenshot of the Gmail inbox interface. On the left, there is a sidebar with the following menu items: "Compose" (highlighted in blue), "Inbox" (highlighted in light blue), "Snoozed", "Sent", "Drafts", "All Mail", and "Spam".





NVIDIA®

Progressive Growing of GANs for Improved
Quality, Stability, and Variation
Karras et al., ICLR'18



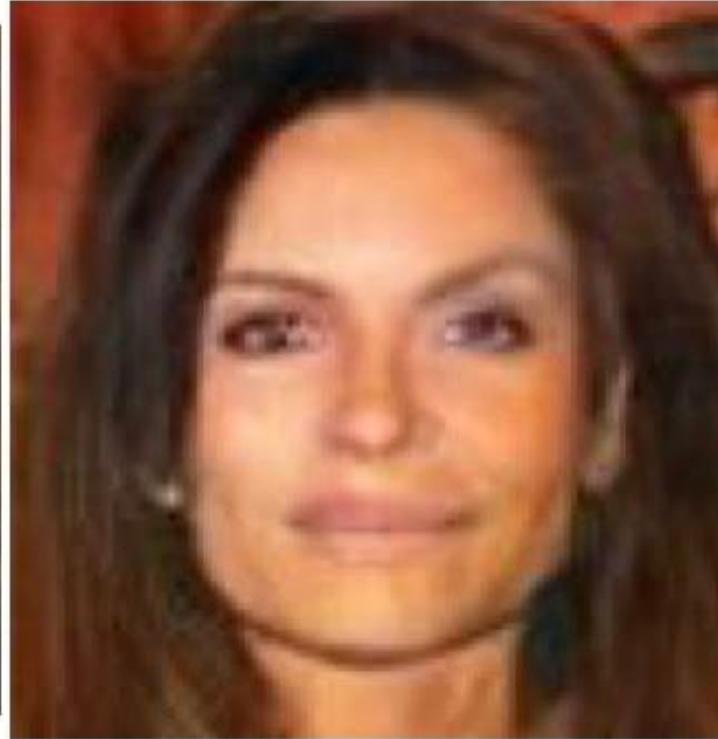
Progressive Growing of GANs for Improved
Quality, Stability, and Variation
Karras et al., ICLR'18



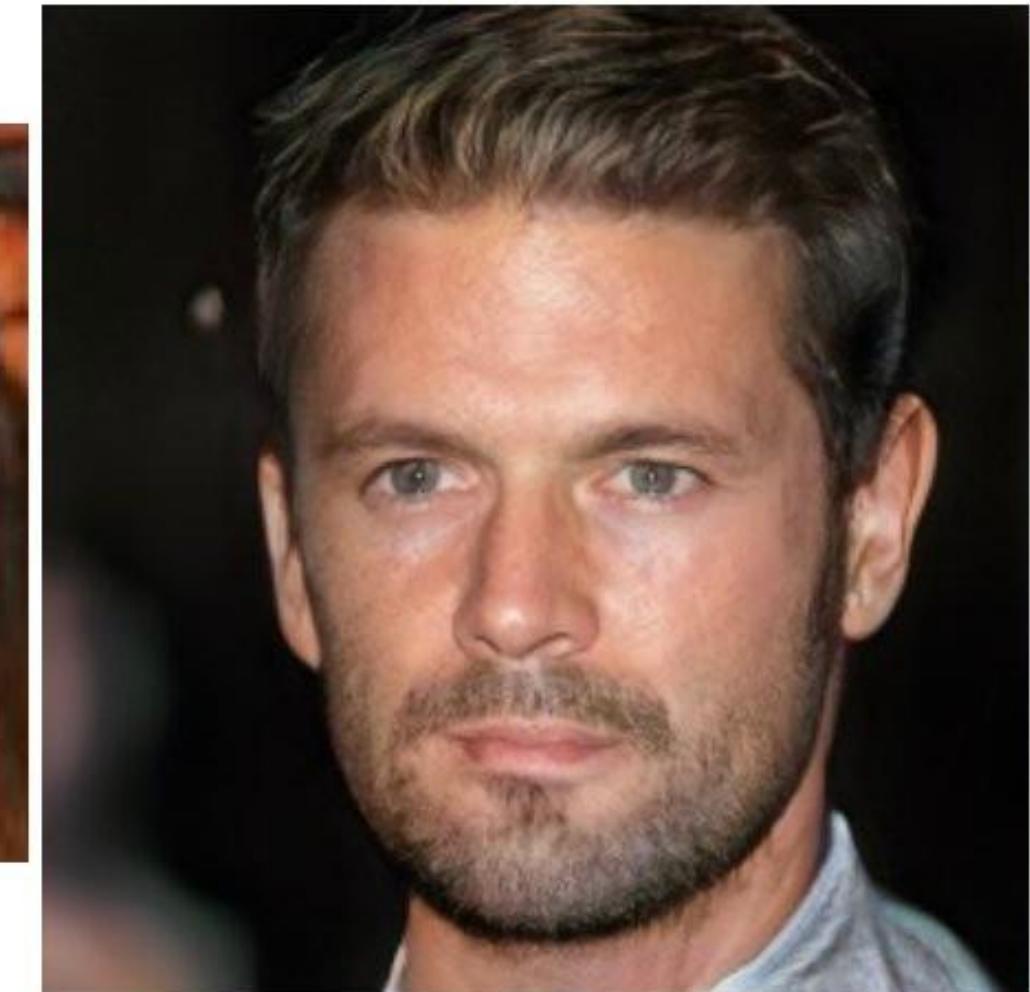
2014



2015



2016



2017



2018



"man in black shirt is playing guitar."



"construction worker in orange safety vest is working on road."



"two young girls are playing with lego toy."



"boy is doing backflip on wakeboard."



"girl in pink dress is jumping in air."



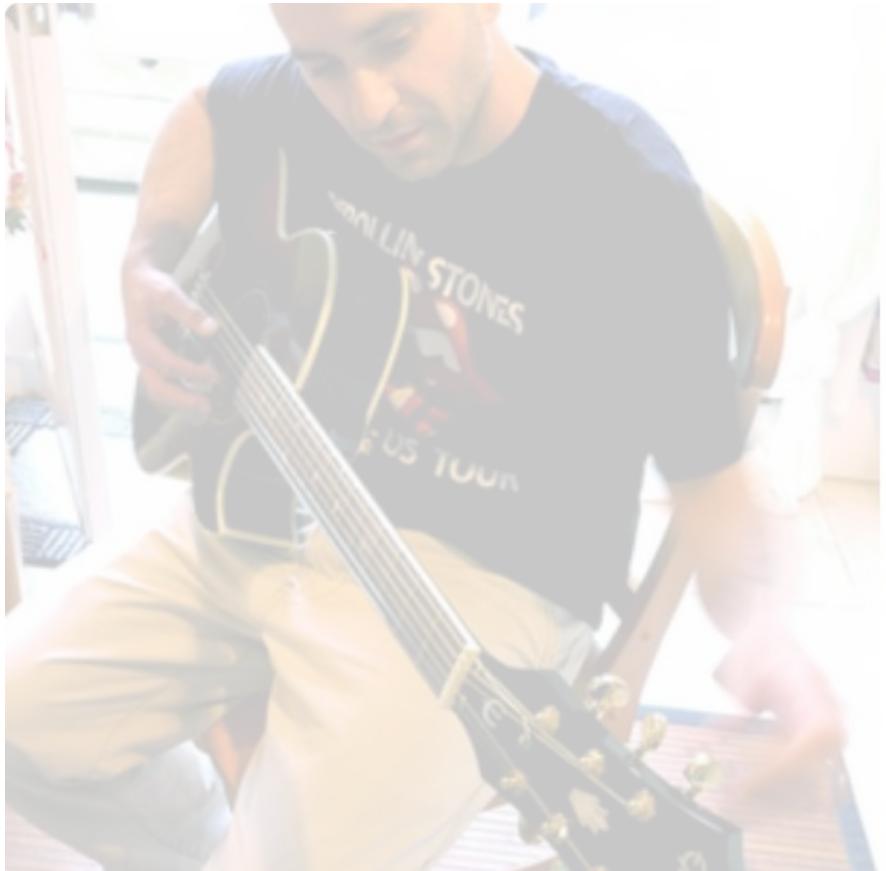
"black and white dog jumps over bar."



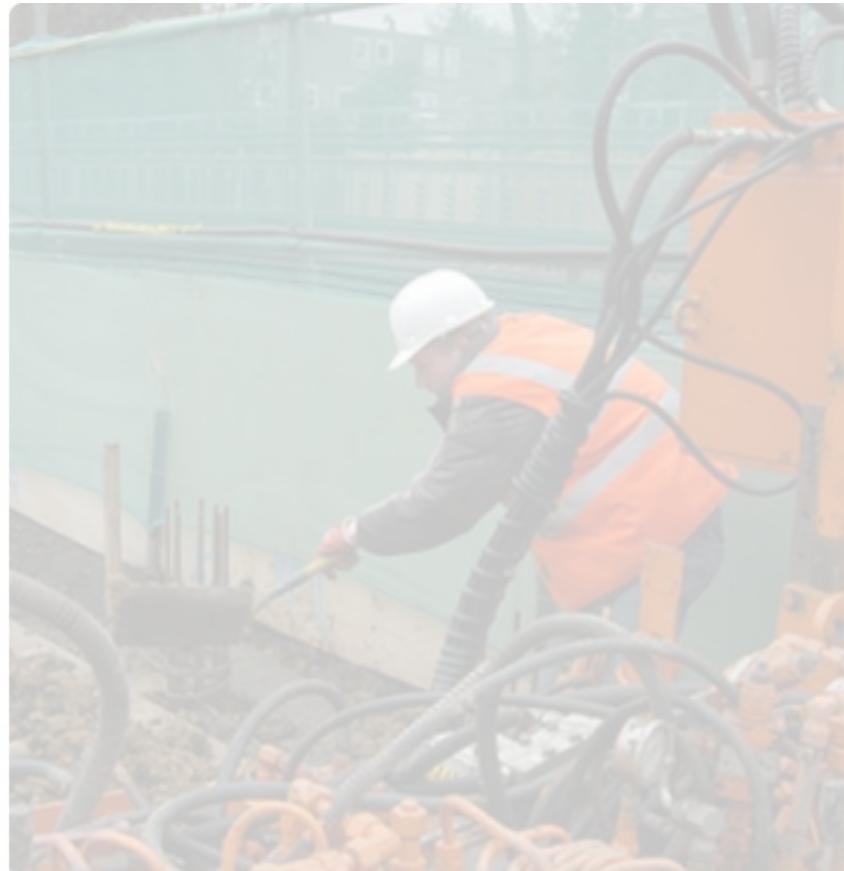
"young girl in pink shirt is swinging on swing."



"man in blue wetsuit is surfing on wave."



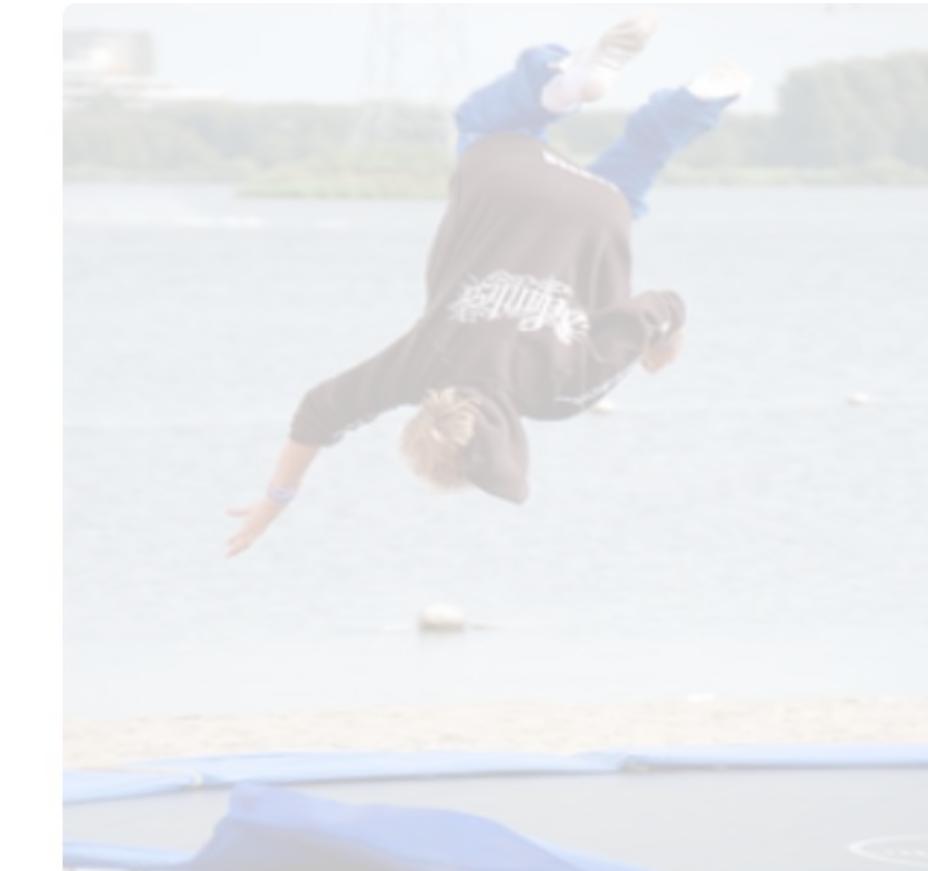
"man in black shirt is playing guitar."



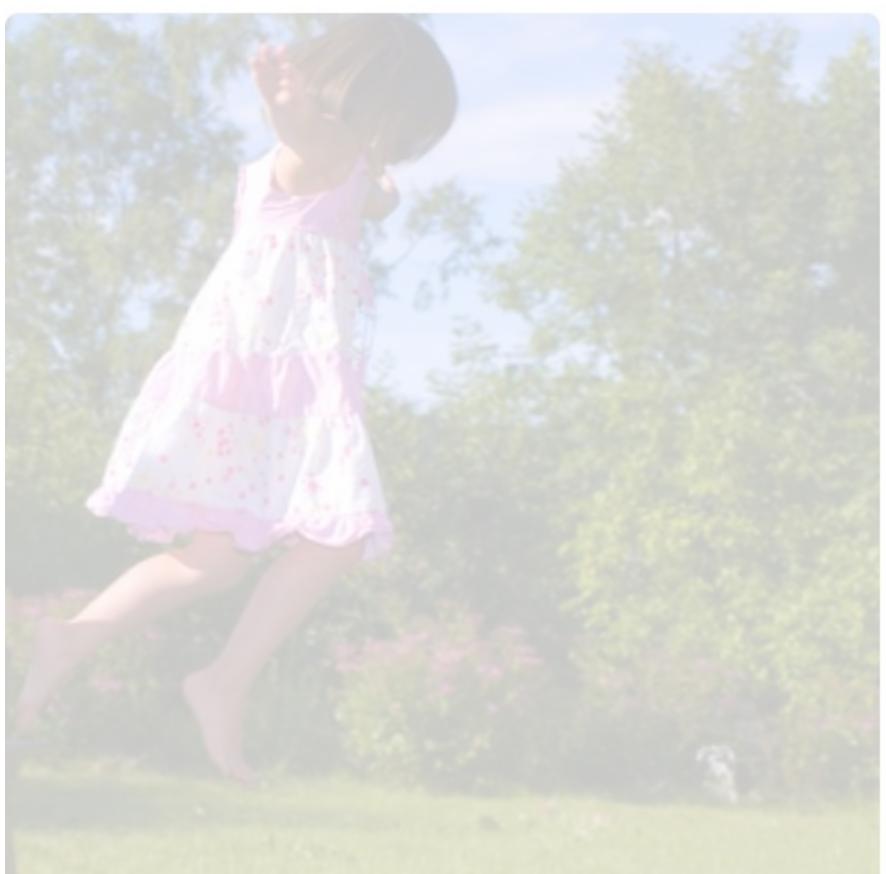
"construction worker in orange safety vest is working on road."



"two young girls are playing with lego toy."



"boy is doing backflip on wakeboard."



"girl in pink dress is jumping in air."



"black and white dog jumps over bar."



"young girl in pink shirt is swinging on swing."



"man in blue wetsuit is surfing on wave."

TEXT PROMPT

an armchair in the shape of an avocado....

AI-GENERATED
IMAGES



Dall-e 2



Stable Diffusion 1.5



<https://huggingface.co/spaces/runwayml/stable-diffusion-v1-5>

<https://openai.com/blog/dall-e/>

Generate an image from a text prompt

TEXT PROMPT

A bowl of soup that looks like a monster spray-painted on the wall

AI-GENERATED
IMAGES



Generate an image from a text prompt

TEXT PROMPT

A bowl of soup that looks like a monster [made out of plasticine](#)

AI-GENERATED
IMAGES



Stable Diffusion 1.5



Generate an image from a text prompt

TEXT PROMPT Two women shaking hands

AI-GENERATED
IMAGES



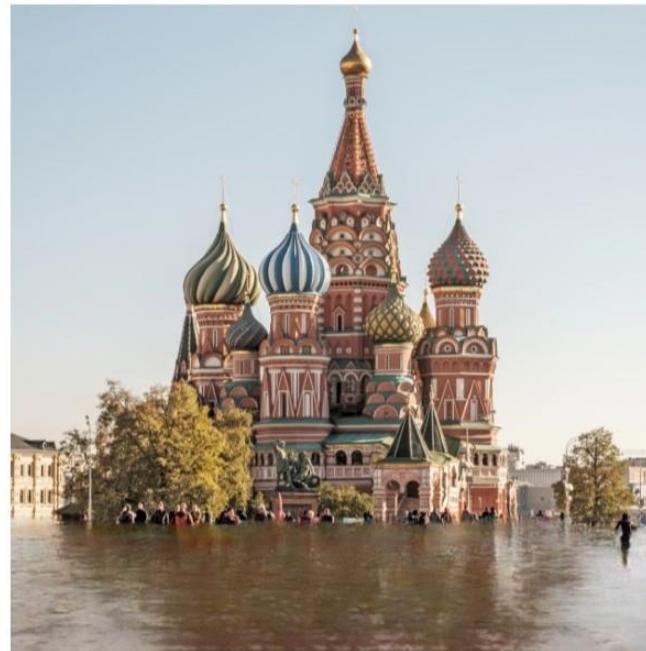
Stable Diffusion 1.5





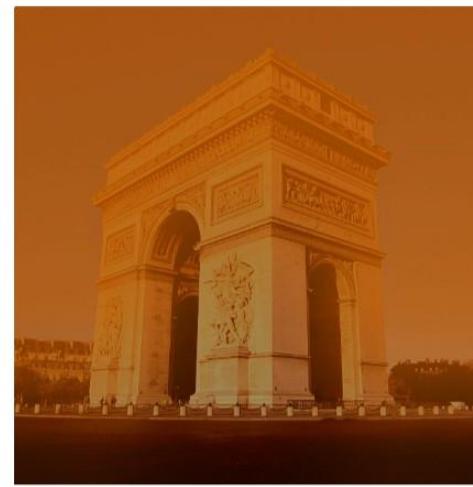
This Climate Does Not Exist

These images
are generated by
artificial intelligence (AI)



Visualize climate change

What now



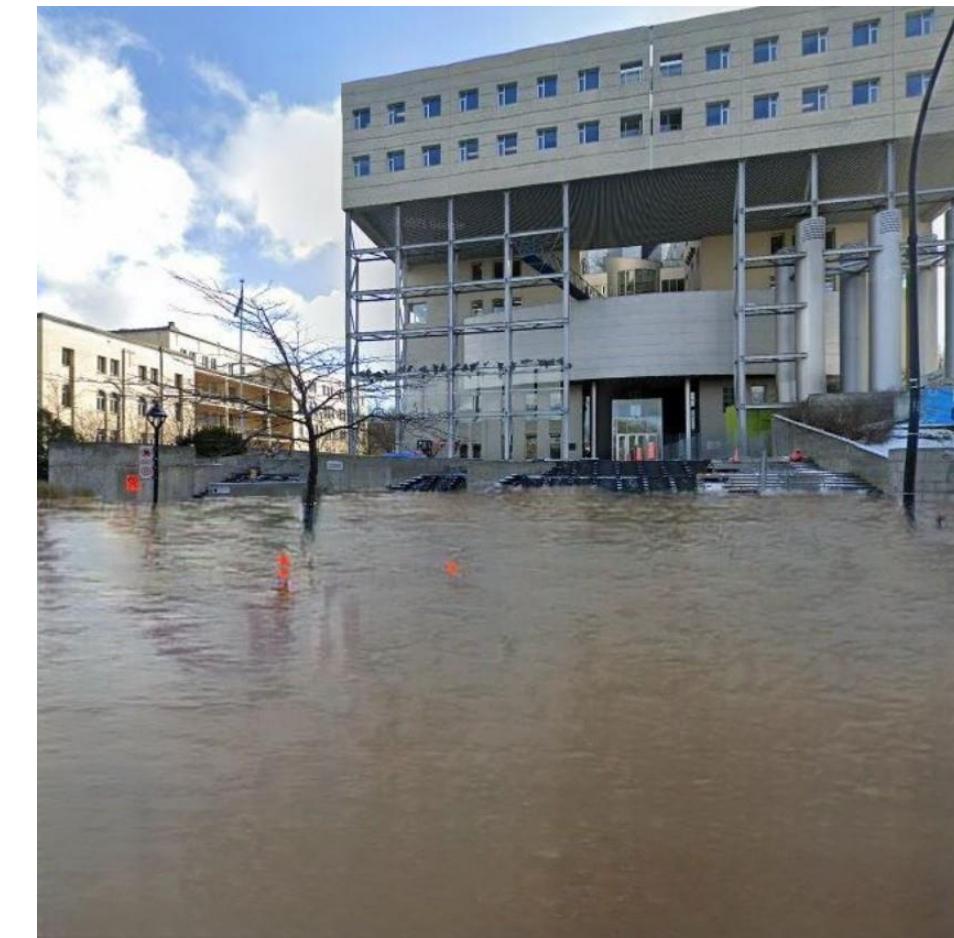
But the environmental
disasters they portray
are very real



Can you imagine these kinds of disasters
happening in your own backyard?



This Climate Does Not Exist



“Lombard street ...”

“view from bottom”



“view from top”



“bird’s eye view”



“in heavy rain”



“Lombard street ...”

“view from bottom”



“view from top”



“bird’s eye view”



“in heavy rain”

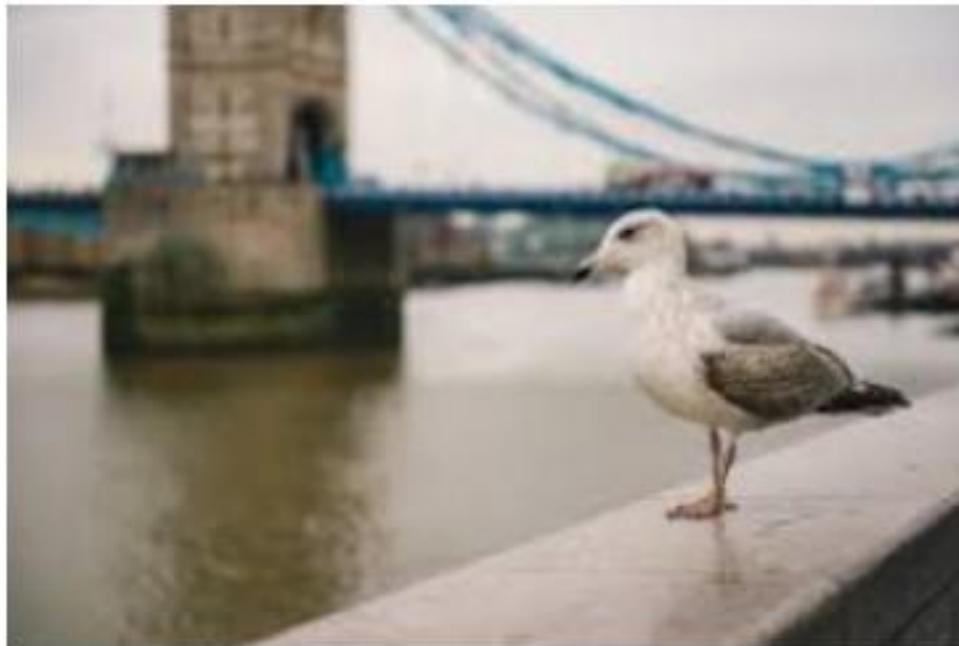


“seagull in front of ...”

“Golden Gate Bridge”



“London Tower Bridge”



“Sydney Harbour Bridge”



“Rialto Bridge”









Robotic observatory makes
fast work of astronomy p. 479

A wet route to
methanol p. 523

Human noise plagues
protected areas p. 531

Science

\$15
5 MAY 2017
sciencemag.org

AAAS

A
♠

DIGITAL CARDS WHIZ

AI beats humans at
challenging poker variant
p. 508



A
♠

AI for video games



Mnih et al.
Nature
Volume 518,
pages 529–533
(26 February 2015)

- **Medicine:** personalized, automate diagnostics
- **Social sciences:** prediction problem (e.g., predict recidivism)
- **Engineering:** to propose new design, evaluate without building
- **Finance:** capture uncertainty, short-term trading
- **Marketing:** to understand and quantify user experience, advertising efficacy
- **Many others:** conservation, social projects, climate change
- **Your domain of expertise...**

Course Introduction & Goals

Logistics

- Course syllabus: <https://davoodwadi.github.io/sessions.html>

Flipped Classroom

- Every week:
 1. Class preparation (Offline):
 - Weekly material (~90 minutes)
 - Reading, watching capsules
 2. Class time (Online):
 - Summary, Q&A, problem solving (120 minutes)

Suggestions for navigating a flipped classroom

- In class: Come prepared
 - Watch the capsules ahead of time
 - Do the readings
 - Write down your questions
- Capsules: Stay active while watching the capsules (e.g., take notes, pause, go back, think of how it fits in the broader context)

Fit with other courses

- HEC
 - PhD level (originally)
 - Computationally oriented
 - Prequel to
 - Machine Learning II: Deep Learning (MATH 60630A)
 - Trustworthy Machine Learning (MATH 80630)
 - Other ML courses in Montreal (U.Montreal, Polytechnique, McGill)
 - More applied (similar to COMP-551@McGill)

Notes on LLMs' use

- **Use for learning code**
- **Not for writing code or text**

An experiment on ChatGPT

Default (GPT-3.5)



I'm doing a homework with the following question:
"Train a neural network using Keras. You will use transfer learning with the ResNet18 backbone and BatchNorm layers between each Conv layer. Find the best learning rate using Adam, AdamW, and SGD with momentum"

An experiment on ChatGPT

Default (GPT-3.5)

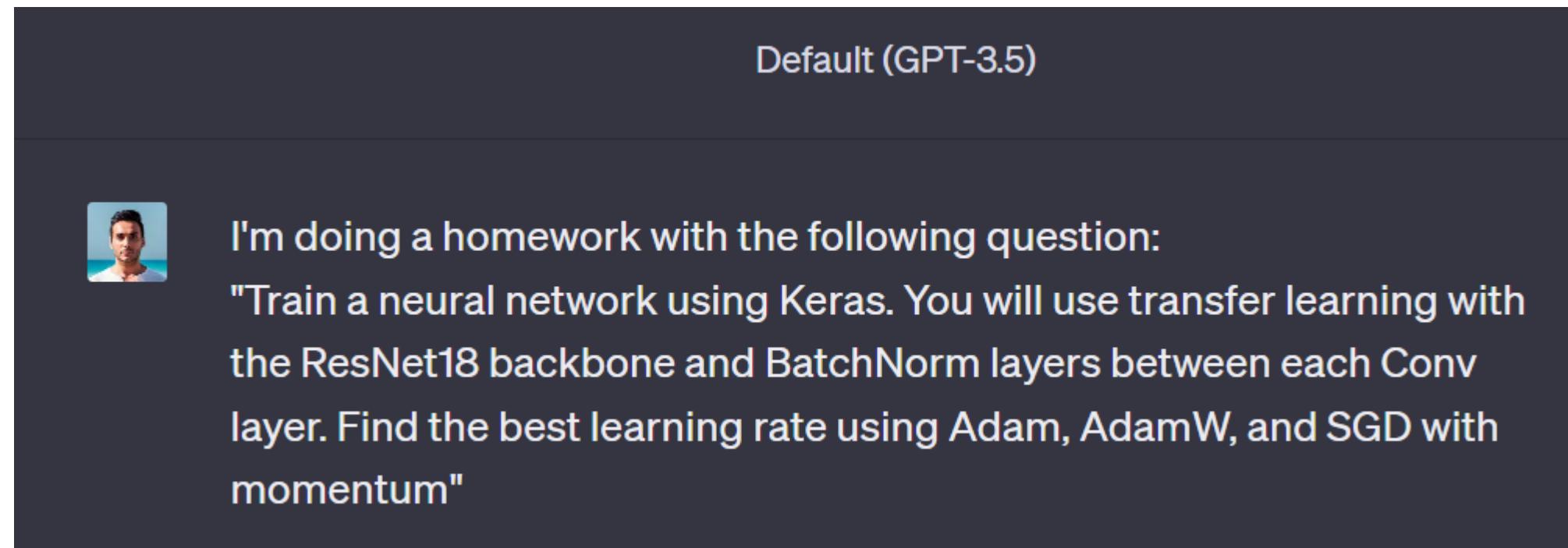


I'm doing a homework with the following question:
"Train a neural network using Keras. You will use transfer learning with the ResNet18 backbone and BatchNorm layers between each Conv layer. Find the best learning rate using Adam, AdamW, and SGD with momentum"

```
import tensorflow as tf
from tensorflow.keras.applications import ResNet18
from tensorflow.keras.layers import Dense, GlobalAveragePooling2D
from tensorflow.keras.models import Model
from tensorflow.keras.optimizers import Adam, AdamW, SGD
from tensorflow.keras.datasets import ...

# Load your dataset here (e.g., CIFAR-10, ImageNet, etc.)
# (X_train, y_train), (X_test, y_test) = ...
```

An experiment on ChatGPT



```
base_model = ResNet18(include_top=False, weights='imagenet', input_shape=(224, 224, 3))

x = base_model.output
x = GlobalAveragePooling2D()(x)
x = Dense(1024, activation='relu')(x)
predictions = Dense(num_classes, activation='softmax')(x) # Adjust num_classes according to your task

model = Model(inputs=base_model.input, outputs=predictions)
```

An experiment on ChatGPT

```
# Learning rate finder (optional)
lr_finder = tf.keras.callbacks.LearningRateScheduler(lambda epoch: 1e-8 * 10**(epoch / 20))

# Define learning rates for experimentation
learning_rates = [1e-4, 1e-3, 1e-2, 1e-1]

for lr in learning_rates:
    # Adam optimizer
    adam_optimizer = Adam(learning_rate=lr)
    model.compile(optimizer=adam_optimizer, loss='categorical_crossentropy', metrics=['accuracy'])
    model.fit(X_train, y_train, batch_size=32, epochs=5, callbacks=[lr_finder])

    # AdamW optimizer
    adamw_optimizer = AdamW(learning_rate=lr)
    model.compile(optimizer=adamw_optimizer, loss='categorical_crossentropy', metrics=['accuracy'])
    model.fit(X_train, y_train, batch_size=32, epochs=5, callbacks=[lr_finder])

    # SGD optimizer with momentum
    sgd_optimizer = SGD(learning_rate=lr, momentum=0.9)
    model.compile(optimizer=sgd_optimizer, loss='categorical_crossentropy', metrics=['accuracy'])
    model.fit(X_train, y_train, batch_size=32, epochs=5, callbacks=[lr_finder])
```

An experiment on ChatGPT

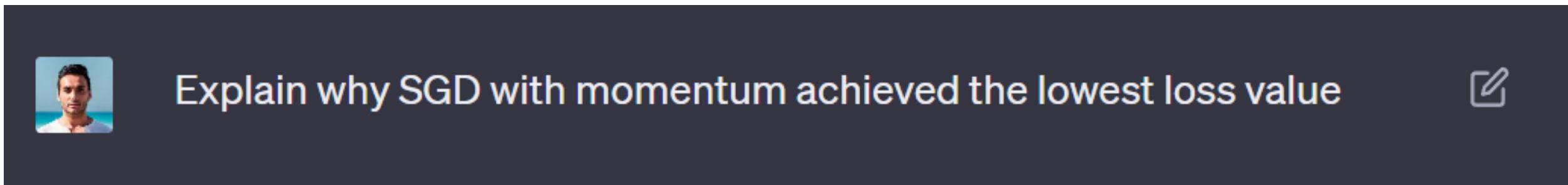
```
import matplotlib.pyplot as plt

# Plot the learning rate vs. loss graph
lrs = learning_rates
losses = [list_of_adam_losses, list_of_adamw_losses, list_of_sgd_losses] # Update with your actual loss values
plt.semilogx(lrs, losses[0], label='Adam')
plt.semilogx(lrs, losses[1], label='AdamW')
plt.semilogx(lrs, losses[2], label='SGD')
plt.xlabel('Learning Rate')
plt.ylabel('Loss')
plt.legend()
plt.show()
```

Visible bugs + hidden bugs

An experiment on ChatGPT

- Use for learning code
- Not for writing code or text



An experiment on ChatGPT

SGD with momentum is an optimization algorithm that builds upon the basic Stochastic Gradient Descent (SGD) method by introducing a momentum term. This momentum term accelerates the optimization process and helps the model overcome the limitations of vanilla SGD, especially in scenarios where the loss surface is rugged or contains many local minima. The achieved lower loss value for SGD with momentum compared to Adam and AdamW could be attributed to several factors:

An experiment on ChatGPT

1. Acceleration of Convergence: The momentum term in SGD introduces a moving average of past gradients. This has the effect of accelerating the convergence towards the minimum, especially in the presence of oscillations or noisy gradients. This acceleration can help the optimizer escape shallow local minima and make quicker progress towards the optimal solution.

2. Inertia: The momentum term allows the optimization process to "carry over" some of the direction from previous steps. This inertia helps the optimization process to navigate through flatter regions of the loss surface more efficiently, potentially leading to faster convergence towards a lower loss.

3. Less Frequent Updates: Adam and AdamW algorithms use more complex mechanisms involving adaptive learning rates and second-moment estimates (RMSprop), which can lead to more frequent updates of the model parameters. In contrast, SGD with momentum can lead to less frequent updates, which may provide the optimization process with a better chance to settle into regions with lower loss values.

4. Applicability to Data: Depending on the nature of the dataset and the specific task, the way gradients are averaged and the learning rate is adjusted can impact the optimization process. SGD with momentum might be better suited to certain datasets and tasks, allowing it to navigate the loss landscape more effectively.

5. Hyperparameter Tuning: The performance of optimization algorithms often depends on hyperparameter settings, such as learning rates, momentum values, and batch sizes. It's possible that for this specific task and dataset, the hyperparameters used for SGD with momentum were more suitable and well-tuned, leading to better convergence.

Notes on LLMs' use

- **HEC Montreal's policy on LLMs:**
<https://www.hec.ca/en/integrity/good-practices-and-ressources/index.html#explo-198752-480086>

If you use content generated by an artificial intelligence tool for an evaluation, you must clearly disclose this by citing your source. Subject to future APA guidance, here is how to quote when using an artificial intelligence tool:

- Template
 - Name of tool used. (YYYY) ChatGPT [version of [day month]] [search question/topic used]. [URL]
- Result
 - OpenAI (2023) ChatGPT (May 27 version) [[How to apply the components of Porter's model to the lithium battery industry](#)].

Notes on LLMs' use

- ChatGPT can be detected using plagiarism software
- ChatGPT is generally
 1. more verbose than a human
 2. has a unique identifiable style



Send a short letter to a finance company saying I'm interested in their [Senior Financial Analyst role](#) 

Notes on LLMs' use

- ChatGPT can be detected using plagiarism software
- ChatGPT is generally
 1. more verbose than a human
 2. has a unique identifiable style

Dear [Company Name] Hiring Team,

I hope this message finds you well. I am writing to express my strong interest in the Senior Financial Analyst role at [Company Name], as advertised on your website.

With a solid educational background in finance and [X years] of experience in financial analysis and reporting, I am excited about the opportunity to contribute my skills and expertise to your **esteemed** organization. I have a **proven track record** of leveraging financial data to drive strategic decision-making, streamline processes, and optimize resource allocation.

From my research into [Company Name]'s reputation for innovation and excellence in the finance industry, I am **particularly drawn** to your commitment to leveraging data-driven insights to guide financial strategies. I am confident that my proficiency in financial modeling, forecasting, and risk assessment aligns well with the demands of the Senior Financial Analyst role.

I am impressed by [Company Name]'s **dedication to fostering** a collaborative and dynamic work environment, where professionals are empowered to make meaningful contributions and grow within the company. The opportunity to work alongside a team of experienced financial experts while utilizing cutting-edge tools and technologies greatly appeals to me...

Short review of linear algebra, statistics, and probabilities

- Based on chapters 2 and 3 of “Deep Learning”

<http://www.deeplearningbook.org/>

Linear algebra

- **Scalar: a single value.**

$$\mathbf{a} \in \mathbb{R}, \mathbf{a} \in \mathbb{N} \quad \mathbf{a} = 3$$

- **Vector: an array of values.**

$$\mathbf{a} \in \mathbb{R}^D, \mathbf{a} \in \mathbb{N}^D \quad \mathbf{a} = \begin{bmatrix} 3 \\ 4 \\ 2 \end{bmatrix}$$

- **Matrix: a table of values.**

$$\mathbf{A} \in \mathbb{R}^{D_1 \times D_2}, \mathbf{A} \in \mathbb{N}^{D_1 \times D_2} \quad \mathbf{A} = \begin{bmatrix} 3 & 4 & 2 \\ 1 & 2 & 9 \end{bmatrix}$$

Indexing notation

- Indexing elements of a vector: a_i

$$a = \begin{bmatrix} 3 \\ 4 \\ 2 \end{bmatrix} \leftarrow a_1$$

Convention:
The first element
is the zero'th.

- Indexing elements of a matrix: a_{ij}

$$A = \begin{bmatrix} 3 & 4 & 2 \\ 1 & 2 & 9 \end{bmatrix}$$

\uparrow
 a_{12}

Simple operations

- Transpose

$$a = \begin{bmatrix} a_0 \\ a_1 \\ a_2 \end{bmatrix} \quad \left| \quad (A_{ij})^\top = A_{ji}$$
$$a^\top = [a_0 \ a_1 \ a_2]$$

- Addition

- Vectors and matrices w. the same shape

$$b = \begin{bmatrix} b_0 \\ b_1 \\ b_2 \end{bmatrix} \quad \left| \quad (A + B)_{ij} = A_{ij} + B_{ij}$$
$$a + b = \begin{bmatrix} a_0 + b_0 \\ a_1 + b_1 \\ a_2 + b_2 \end{bmatrix}$$

Simple operations

- Multiply by a scalar

$$\alpha \mathbf{a} = \begin{bmatrix} \alpha \mathbf{a}_0 \\ \alpha \mathbf{a}_1 \\ \alpha \mathbf{a}_2 \end{bmatrix}$$

- Vector product.

- The dot product

$$\mathbf{a}^\top \mathbf{a} = \sum_i \mathbf{a}_i \mathbf{a}_i$$

- Note: it yields a scalar.

- Element-wise product:

$$\mathbf{a} \odot \mathbf{a} = \begin{bmatrix} \mathbf{a}_0 \mathbf{a}_0 \\ \mathbf{a}_1 \mathbf{a}_1 \\ \mathbf{a}_2 \mathbf{a}_2 \end{bmatrix}$$

- Also known as Hadamard product

Operations

- Matrix product (dot product):

$$C_{ij} = \sum_k A_{ik} B_{kj}$$

- A's columns must equal B's rows (order is important)

$$\mathbf{A} \in \mathbb{R}^{D_1 \times D_2}, \mathbf{B} \in \mathbb{R}^{D_2 \times D_3}$$

- Distributive: $\mathbf{A}(\mathbf{B} + \mathbf{C}) = \mathbf{AB} + \mathbf{AC}$
- Associative: $\mathbf{A}(\mathbf{BC}) = (\mathbf{AB})\mathbf{C}$
- Product of transpose: $(\mathbf{AB})^\top = \mathbf{B}^\top \mathbf{A}^\top$

Inverse

- We denote a matrix's inverse as A^{-1}
- A matrix has an inverse iff:
 - it's square. $D_1 = D_2$
 - its columns are linearly independent.
 - No column can be recovered using a combination of other columns
- Inverses are useful to solve systems of equations:

$$Ax = b \quad x = A^{-1}b$$

A square matrix
not invertible is *singular*

Norms

- L^p norm. Size of a vector (or matrix)

$$\| \mathbf{a} \|_p = \left(\sum_i |a_i|^p \right)^{1/p}$$

- Standard norms in ML:

- Euclidean norm ($p=2$)

$$\| \mathbf{a} \|_2 = \sqrt{\left(\sum_i |a_i|^2 \right)}$$

- Dot product w. 2-norm: $\mathbf{a}^\top \mathbf{b} = \|\mathbf{a}\|_2 \|\mathbf{b}\|_2 \cos \theta_{ab}$

- Frobenius norm (matrix): $\| \mathbf{A} \|_2 = \sqrt{\left(\sum_i \sum_j |a_{ij}|^2 \right)}$

Special matrices & vectors

- Identity. Denoted I_n .

$$I_3 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

- All zeros except for ones on the main diagonal.
- Symmetric: $A = A^\top$
- Unit vector: $\|a\|_2 = 1$
- Orthogonal vectors: $a^\top b = 0$
- Orthonormal vectors: unit and orthogonal $A^\top A = AA^\top = I$
- Orthogonal matrix: Orthonormal rows & columns

- Skip eigendecomposition, SVD, pseudo-Inverse, determinants (Sections 2.7–2.11).
 - We will get back to them if/when needed in the course.

- On to probabilities
- Chapter 3 of “Deep Learning”
 - I’ve adapted some of the lecture slides from the book.
 - Thanks to Ian Goodfellow for providing slides.

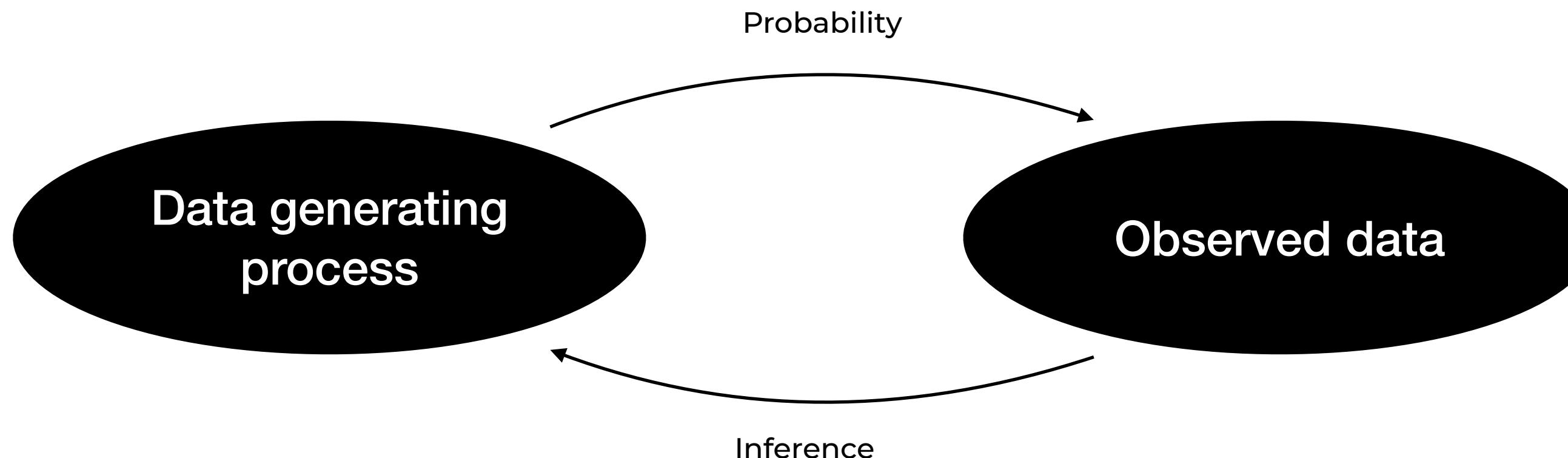
Why probabilities?

- To capture uncertainty

E.g., What time will I get home tonight?

- Probabilities provide a formalism for making statements about “data generating processes” (L. Wasserman)

E.g., what happens when I flip a fair coin?



The example

- Generate data by throwing a fair die.
- What do we know about a single throw?
 - 6 possible outcomes. (**sample space**)
 - Each outcome (e.g., 1). (**element, state**)
 - A subset of outcomes (e.g., <3). (**event**)
 - Outcomes are equiprobable. (**uniform distribution**)

Random variables and probabilities

- A random variable (r.v.) is a probabilistic outcome.
 - For example,
 - Die throw (X)
 - The actual outcome is $\in \{1, 2, 3, 4, 5, 6\}$. (x)
 - A probability function (P) assigns a real number to each possible event: $P(x) \geq 0, \forall x \in X$

$$P(\bigcup x) = 1$$

Discrete RVs

- An RV is discrete if it takes a finite number of values¹

$$\begin{aligned} P(x = x_i) &\geq 0, \forall i \\ \sum_i P(x = x_i) &= 1 \end{aligned}$$

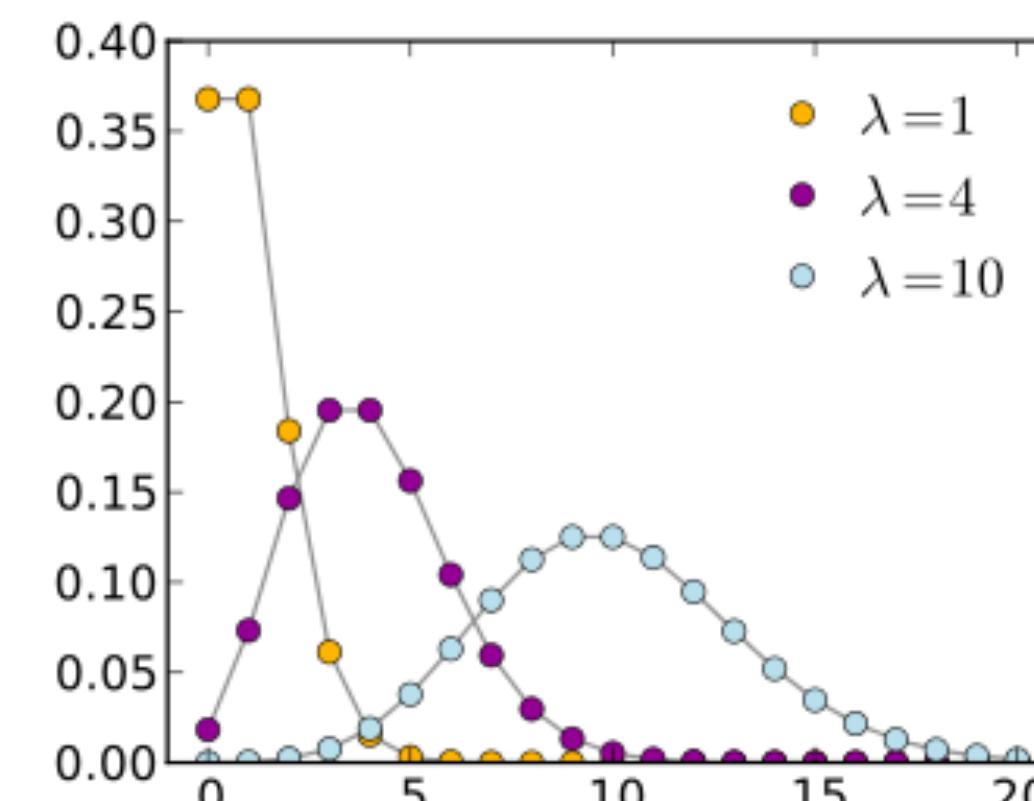
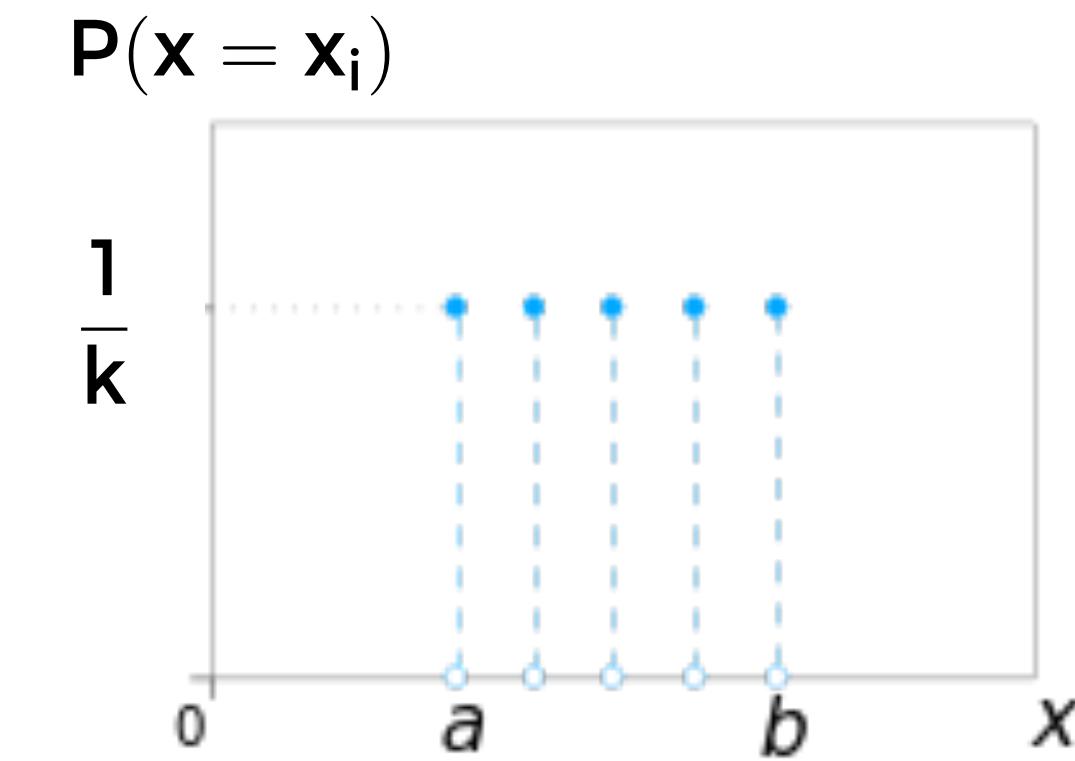
- E.g., uniform distribution:

$$P(x = x_i) = \frac{1}{k}, \forall i$$

- E.g., Poisson distribution:

$$P(x = x_i; \lambda) = \frac{\lambda^{x_i} e^{-\lambda}}{x_i!}$$

1. technically: it must be countable



Continuous RVs

- An RV is continuous if $f(x) \geq 0, \forall x \in X$

$$\int_{-\infty}^{\infty} f(x)dx = 1$$

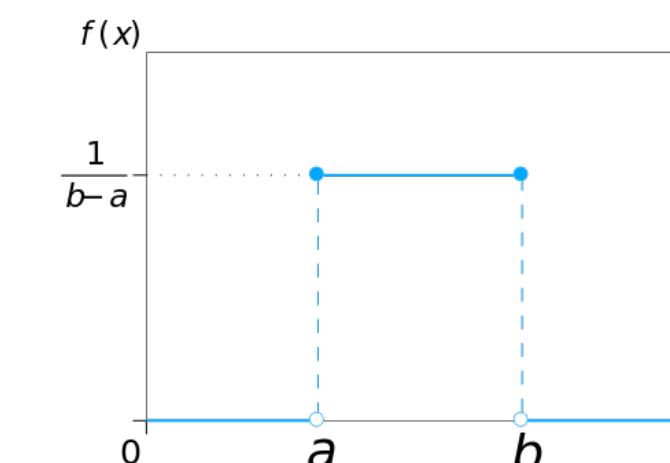
$$P(a < x < b) = \int_a^b f(x)dx$$

- $f(x)$ is a probability density function (PDF)

- E.g., (continuous) uniform distribution:

$$u(x; a, b) = \begin{cases} \frac{1}{b-a} & \text{if } x \in [a, b] \\ 0 & \text{otherwise} \end{cases}$$

- E.g., Gaussian distribution



from: wikipedia.org

A few useful properties

(shown for discrete variables for simplicity)

- **Sum rule:** $P(X) = \sum_Y P(X, Y)$
- **Product rule:** $P(X, Y) = P(X | Y)P(Y)$
- **Chain rule:** $P(X_1, \dots, X_n) = \prod_{i=1}^n P(X_i | X_1, \dots, X_{i-1})P(X_1)$
- **If x and y are independent:** $P(X, Y) = P(X)P(Y)$
- **Bayes' Rule:** $P(Y | X) = \frac{P(X | Y)P(Y)}{P(X)}$

Moments

- **Expectation:** $\mathbb{E}[X] = \sum_i P(x = x_i)x_i \quad \mathbb{E}[aX] = a\mathbb{E}[X]$
- **Variance:** $\sigma^2 = \mathbb{E}[(X - \mathbb{E}[X])^2]$
- **Covariance:** $\text{Cov}(X, Y) = \mathbb{E} [(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])]$
- **correlation:** $\rho(x, y) = \frac{\text{Cov}(X, Y)}{\sigma_x \sigma_y}$

Further Reading

- Prologue to “The Master Algorithm”
<http://homes.cs.washington.edu/~pedrod/Prologue.pdf>
- Ch. 1 of Hastie et al.
- Math Preparation
 - Ch.2 of Pattern Recognition and Machine Learning [PRML]
 - Ch.2-3 of Deep Learning [DL]
 - Slightly more advanced:

<http://www.cs.mcgill.ca/~dprecup/courses/ML/Materials/prob-review.pdf>

<http://www.cs.mcgill.ca/~dprecup/courses/ML/Materials/linalg-review.pdf>