

When Models Meet Data

Jo Ciucă

Australian National University

comp36706670@anu.edu.au

Class reps



Judy Xie

Jinghan Wu

Qianhui Zhang

Richard Susilo

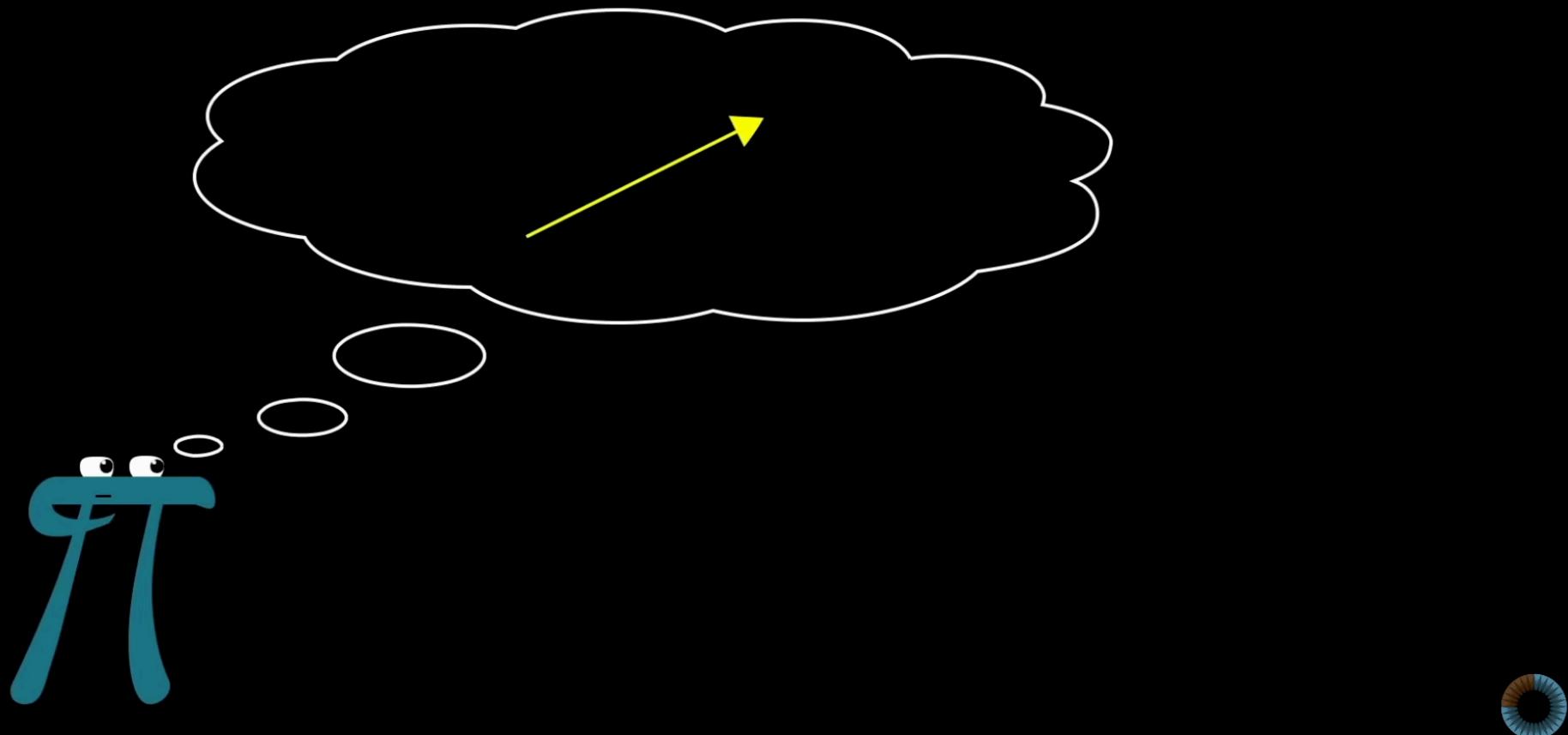
Alvaro Bhirawa

Arthur Zhao

Outline

- Gram-Schmidt Orthogonalization
- Models learning from data
- Model meets data: finding parameters
- Empirical Risk Minimisation

Think of individual vectors as arrows



Credit: 3blue1brown

3.8.3 Gram-Schmidt Orthogonalization

- Consider a basis of \mathbb{R}^2

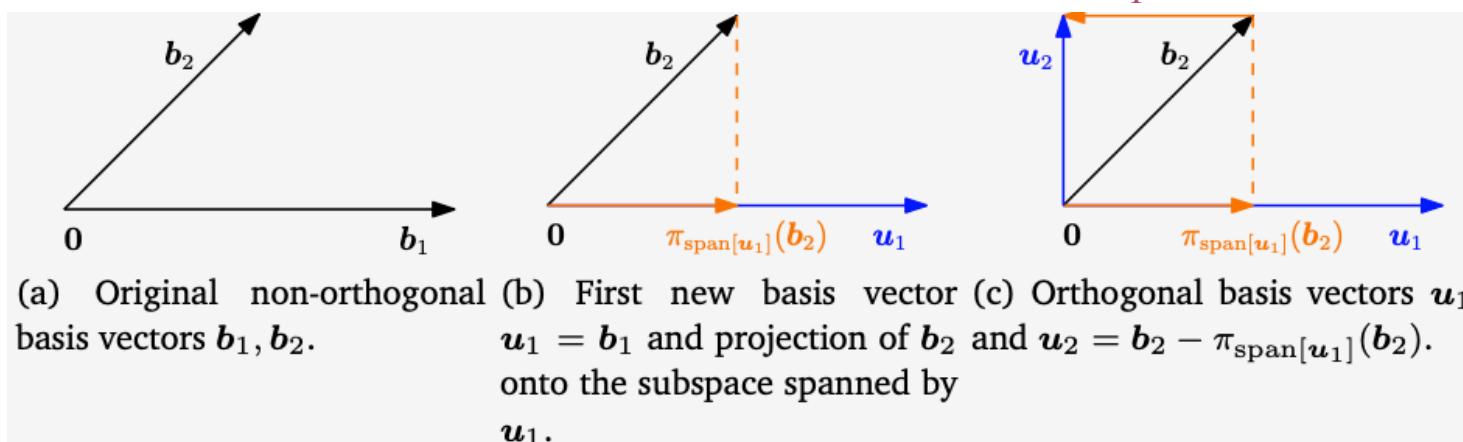
$$\mathbf{b}_1 = \begin{bmatrix} 2 \\ 0 \end{bmatrix} \quad \mathbf{b}_2 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

- Using the Gram-Schmidt method, we construct an **orthogonal** basis $(\mathbf{u}_1, \mathbf{u}_2)$ of \mathbb{R}^2 as follows (using dot product).

$$\mathbf{u}_1 := \mathbf{b}_1 = \begin{bmatrix} 2 \\ 0 \end{bmatrix}$$

$$\mathbf{u}_2 := \mathbf{b}_2 - \pi_{\text{span}[\mathbf{u}_1]}(\mathbf{b}_2) = \mathbf{b}_2 - \frac{\mathbf{u}_1 \mathbf{u}_1^T}{\|\mathbf{u}_1\|^2} \mathbf{b}_2 = \begin{bmatrix} 1 \\ 1 \end{bmatrix} - \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

- We immediately see that $\mathbf{u}_1, \mathbf{u}_2$ are orthogonal, i.e., $\mathbf{u}_1^T \mathbf{u}_2 = 0$



3.8.3 Gram-Schmidt Orthogonalization

- Constructively transform basis $(\mathbf{b}_1, \dots, \mathbf{b}_n)$ of an n -dim vector space V into an orthogonal/orthonormal basis $(\mathbf{u}_1, \dots, \mathbf{u}_n)$ of V .

$$\text{span}[\mathbf{b}_1, \dots, \mathbf{b}_n] = \text{span}[\mathbf{u}_1, \dots, \mathbf{u}_n]$$

- The process iterates as follows

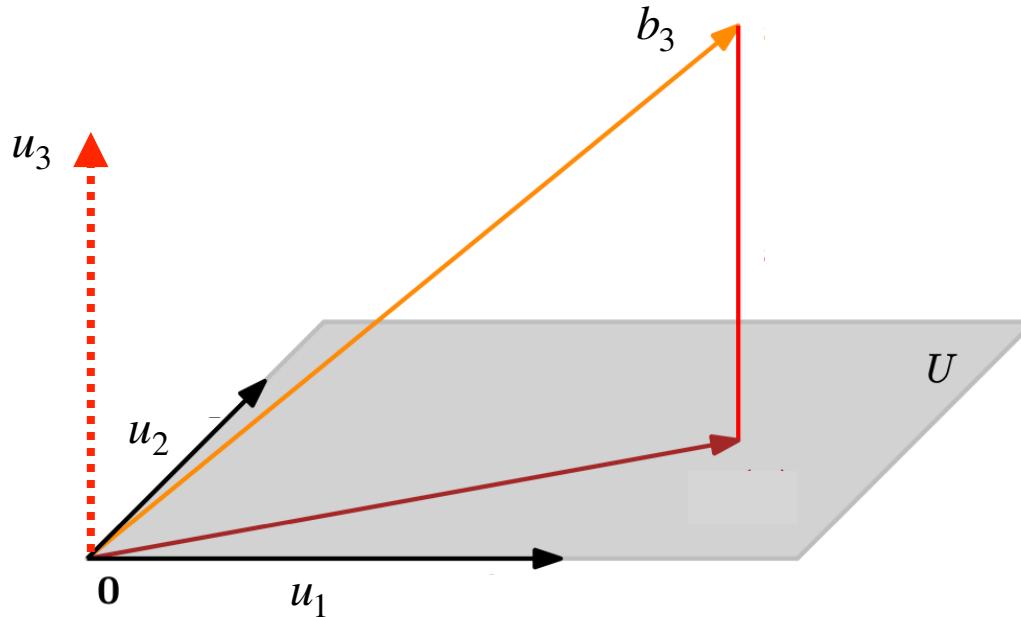
$$\mathbf{u}_1 := \mathbf{b}_1$$

$$\mathbf{u}_k := \mathbf{b}_k - \pi_{\text{span}[\mathbf{u}_1, \dots, \mathbf{u}_{k-1}]}(\mathbf{b}_k), \quad k = 2, \dots, n$$

- The k th basis vector \mathbf{b}_k is projected onto the subspace spanned by the first $k-1$ constructed orthogonal vectors $\mathbf{u}_1, \dots, \mathbf{u}_{k-1}$.
- This projection is then subtracted from \mathbf{b}_k and yields a vector \mathbf{u}_k that is orthogonal to the $(k-1)$ -dim subspace spanned by $\mathbf{u}_1, \dots, \mathbf{u}_{k-1}$.
- If we normalize \mathbf{u}_k , we obtain an ONB where $\|\mathbf{u}_k\| = 1$ for $k = 1, \dots, n$.

3.8.3 Gram-Schmidt Orthogonalization

- Choose one direction as your first $\mathbf{u}_1 := \mathbf{b}_1$
- Project \mathbf{b}_k onto the subspace spanned by the first $k - 1$ constructed orthogonal vectors $\mathbf{u}_1, \dots, \mathbf{u}_{k-1}$.
- Subtract the projection from \mathbf{b}_k and yields a vector \mathbf{u}_k that is orthogonal to the $(k - 1)$ -dim subspace spanned by $\mathbf{u}_1, \dots, \mathbf{u}_{k-1}$
- Normalize \mathbf{u}_k , we obtain an ONB where $\|\mathbf{u}_k\| = 1$ for $k = 1, \dots, n$.





Is there a line
that goes through
the 3 points?

$$y = mx + b$$

let's choose 3 points: $(1, 1), (2, 3), (3, 3)$

$$(1, 1) : m \cdot 1 + b = 1$$

$$(2, 3) : m \cdot 2 + b = 3 \Rightarrow \begin{bmatrix} 1 & 1 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} m \\ b \end{bmatrix} = \begin{bmatrix} 1 \\ 3 \end{bmatrix}$$

$$(3, 3) : m \cdot 3 + b = 3$$

$$\underbrace{\begin{bmatrix} 1 & 1 \\ 2 & 1 \\ 3 & 1 \end{bmatrix}}_A \underbrace{\begin{bmatrix} m \\ b \end{bmatrix}}_u = \begin{bmatrix} 1 \\ 3 \\ 3 \end{bmatrix}$$

But we can solve $A \tilde{u} = p$, where
 p is closest vector to \tilde{u} that lies in
 $C(A)$. And get values for m, b .

Check your understanding

- (A) Orthogonal projections are linear projections.
- (B) The result will no longer change when applying orthogonal projection multiple times (>1).
- (C) Given a subspace to project on, orthogonal projection gives the minimum information loss ($\| \cdot \|_2$).
- (D) Gram-Schmidt Orthogonalization outputs the same number of basis vectors as the input.
- (E) Projections allow us to visualize better and understand high-dimensional data.

Check your understanding

- (A) Orthogonal projections are linear projections. YES
- (B) The result will no longer change when applying orthogonal projection multiple times (>1). YES
- (C) Given a subspace to project on, orthogonal projection gives the minimum information loss (I_2). YES
- (D) Gram-Schmidt Orthogonalization outputs the same number of basis vectors as the input. YES
- (E) Projections allow us to visualise better and understand high-dimensional data. YES

Further resources for Linear Algebra to hone the intuition

Prof Gilbert Strang MIT course on
Linear Algebra

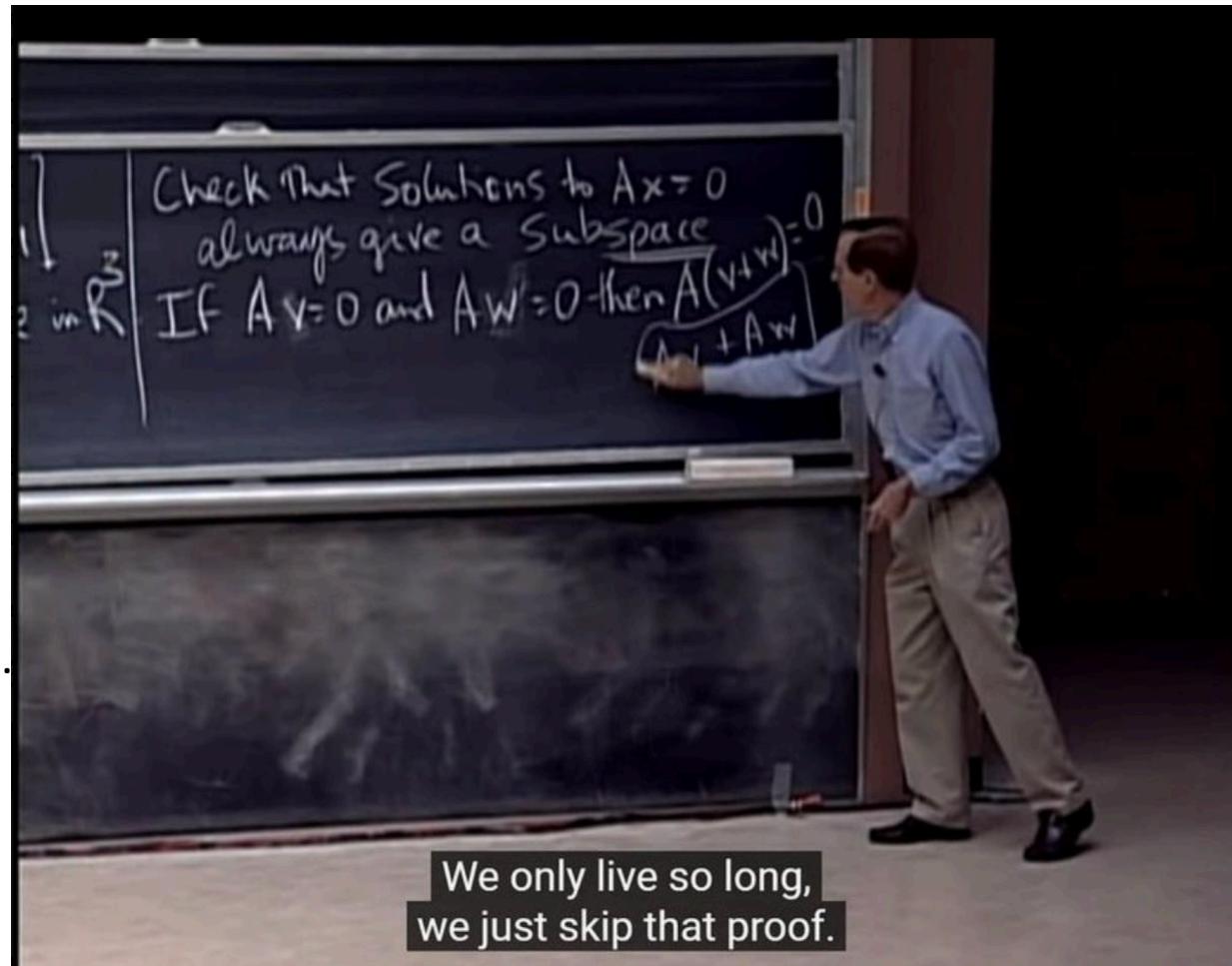
3blue1brown Chapter I

Linear Algebra done right - Axler

Linear Algebra - Serge Lang

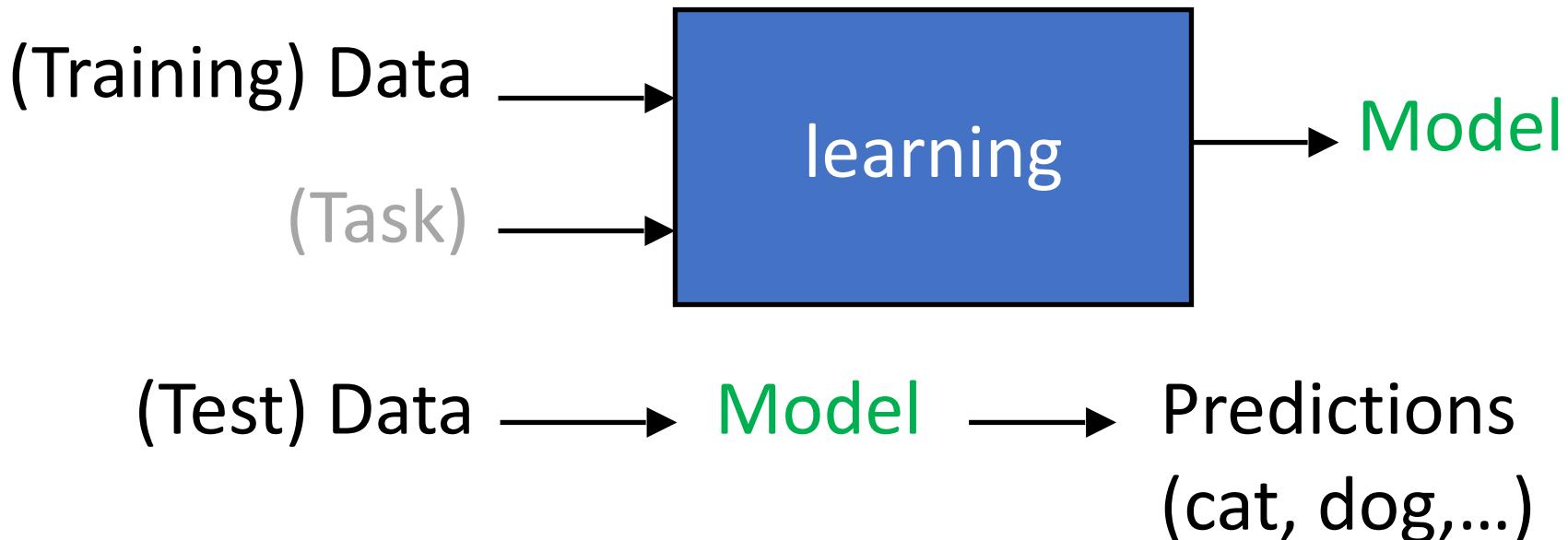
Linear Algebra - Shilov

Deep Learning Book - Bengio et al.



8.1 Data, Models, and Learning

- A machine learning system has three major components:
 - Data, models, learning
-
- A **model** is obtained by **learning** from the training **data**
 - A **prediction** is made by applying a learned **model** on test **data**



8.1 Data, Models, and Learning

- We aim to learn **good** models.
- How is **good** defined? We need to have performance metrics on the test data. Examples include:
 - Classification accuracy
 - Distance from the ground truth
 - Test time (efficiency)
 - Model size
 -
- New performance metrics are constantly being proposed by the machine learning community.

8.1.1 Data as Vectors

- Data, read by computers, should be in a numerical format.
- See the tabular format below

Name	Gender	Degree	Postcode	Age	Annual salary
Aditya	M	MSc	W21BG	36	89563
Bob	M	PhD	EC1A1BA	47	123543
Chloé	F	BEcon	SW1A1BH	26	23989
Daisuke	M	BSc	SE207AT	68	138769
Elisabeth	F	MBA	SE10AA	33	113888

- Row: an instance
- Column: a particular feature
- Apart from tabular format, machine learning can be applied to many types of data, e.g., genomic sequences, text and image contents of a webpage, and social media graphs, citation networks...

- We convert the table into numerical format

Name	Gender	Degree	Postcode	Age	Annual salary
Aditya	M	MSc	W21BG	36	89563
Bob	M	PhD	EC1A1BA	47	123543
Chloé	F	BEcon	SW1A1BH	26	23989
Daisuke	M	BSc	SE207AT	68	138769
Elisabeth	F	MBA	SE10AA	33	113888



	Gender ID	Degree	Latitude (in degrees)	Longitude (in degrees)	Age	Annual Salary (in thousands)
	-1	2	51.5073	0.1290	36	89.563
	-1	3	51.5074	0.1275	47	123.543
	+1	1	51.5071	0.1278	26	23.989
	-1	1	51.5075	0.1281	68	138.769
	+1	2	51.5074	0.1278	33	113.888

- Gender is quantized to -1 and +1
- Degree from BS, MS to PhD: 1, 2, 3
- Postcode corresponds to Latitude and Longitude on the map
- Name is removed because of privacy and because it does not contain useful information for the machine learning system. (exceptions? See [1])

[1] Chen et al., What's in a Name? First Names as Facial Attributes. CVPR 2013

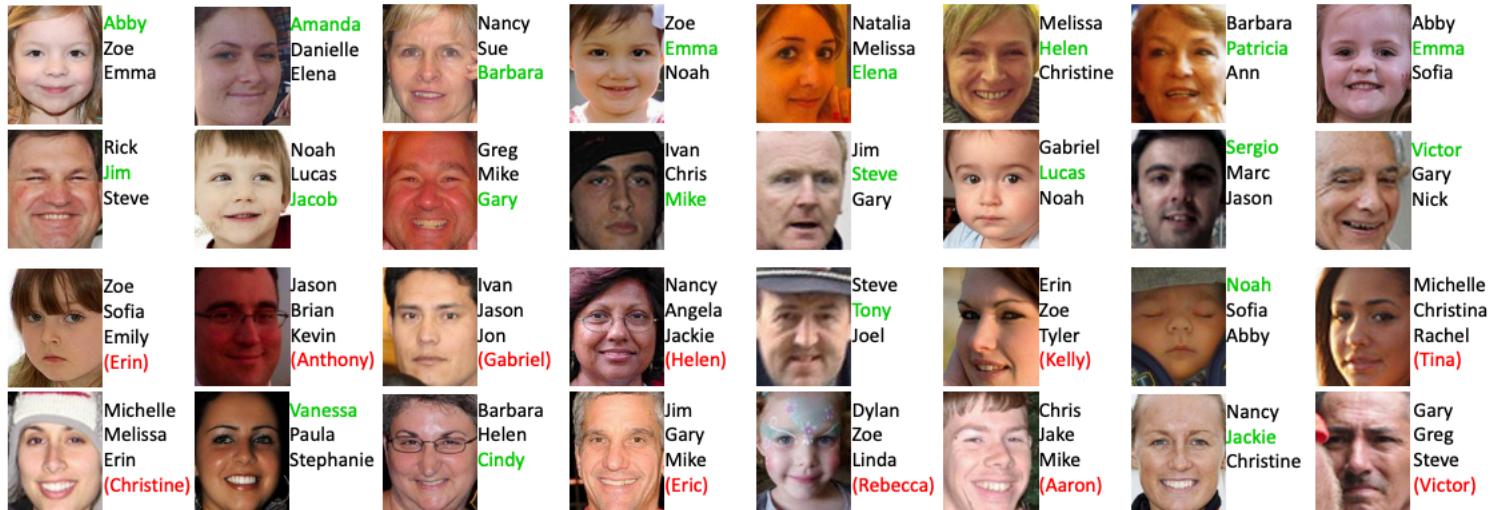


Figure 4: The top-3 predicted names for some face images. The correct prediction is highlighted in green, while the actual first name is shown in red if it is not ranked within the top-3 predictions. The first 2 rows give some good examples where our top-3 predictions include the actual name, and the bottom 2 rows are randomly selected from our test set. Even when our predictions are wrong, reasonable names are predicted (e.g., appropriate gender or age).

- We use N to denote the number of examples in a dataset and index the examples with lowercase $n = 1, \dots, N$

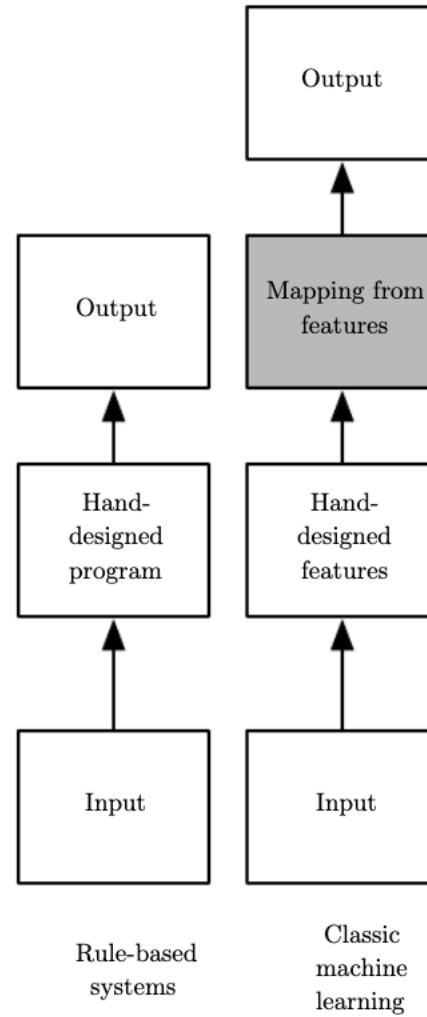
Gender ID	Degree	Latitude (in degrees)	Longitude (in degrees)	Age	Annual Salary (in thousands)
-1	2	51.5073	0.1290	36	89.563
-1	3	51.5074	0.1275	47	123.543
+1	1	51.5071	0.1278	26	23.989
-1	1	51.5075	0.1281	68	138.769
+1	2	51.5074	0.1278	33	113.888

- Each row is a particular individual x_n referred to as an **example** or **data point** in machine learning
- The subscript n refers to the fact that this is the n th example out of a total of N examples in the dataset
- Each column represents a particular feature of interest about the example, and we index the features as $d = 1, \dots, D$
- Each example is a D -dimensional vector

Classical Machine Learning vs Representation Learning

8	9	0	1	2	3	4	7	8	9	0	1	2	3	4	5	6	7	8	6
4	2	6	4	7	5	5	4	7	8	9	2	9	3	9	3	8	2	0	5
0	1	0	4	2	6	5	3	5	3	8	0	0	3	4	1	5	3	0	8
3	0	6	2	7	1	1	8	1	7	1	3	8	9	7	6	7	4	1	6
7	5	1	7	1	9	8	0	6	9	4	9	9	3	7	1	9	2	2	5
3	7	8	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7	0
1	2	3	4	5	6	7	8	9	8	1	0	5	5	1	9	0	4	1	9
3	8	4	7	7	8	5	0	6	5	5	3	3	3	9	8	1	4	0	6
1	0	0	6	2	1	1	3	2	8	8	7	8	4	6	0	2	0	3	6
8	7	1	5	9	9	3	2	4	9	4	6	5	3	2	8	5	9	4	1
6	5	0	1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7
8	9	0	1	2	3	4	5	6	7	8	9	6	4	2	6	4	7	5	5
4	7	8	9	2	9	3	9	3	8	2	0	9	8	0	5	6	0	1	0
4	2	6	5	5	4	3	4	1	5	3	0	8	3	0	6	2	7	1	
1	8	1	7	1	3	8	5	4	2	0	9	7	6	7	4	1	6	8	4
7	5	1	2	6	7	1	9	8	0	6	9	4	9	9	6	2	3	7	1
9	2	2	5	3	7	8	0	1	2	3	4	5	6	7	8	0	1	2	3
4	5	6	7	8	0	1	2	3	4	5	6	7	8	9	2	1	2	1	3
9	9	8	5	3	7	0	7	7	5	7	9	9	4	7	0	3	4	1	4
4	7	5	8	1	4	8	4	1	8	6	4	4	6	3	5	7	2	5	9

MNIST dataset, ~70k images



- Consider the problem of predicting annual salary from age

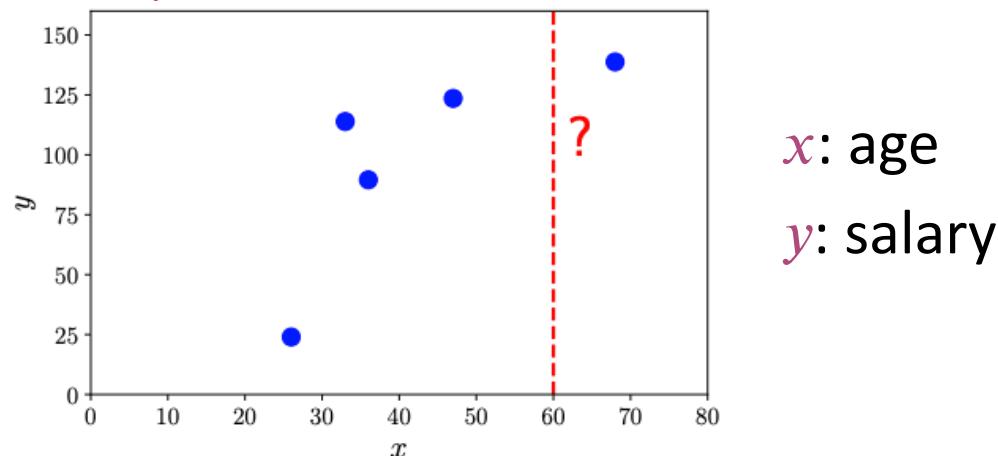
D columns

Gender ID	Degree	Latitude (in degrees)	Longitude (in degrees)	Age	Annual Salary (in thousands)
N rows	-1	2	51.5073	0.1290	36
	-1	3	51.5074	0.1275	47
	+1	1	51.5071	0.1278	26
	-1	1	51.5075	0.1281	68
	+1	2	51.5074	0.1278	33

- A **supervised learning** algorithm
- We have a **label** y_n (the salary) associated with each **example** x_n (e.g., age).
- A dataset is written as a set of example-label pairs
 $\{(x_1, y_1), \dots, (x_n, y_n), \dots, (x_N, y_N)\}$
- The table of examples $\{x_1, \dots, x_N\}$ are concatenated and written as $X \in \mathbb{R}^{N \times D}$

We are interested in:

What is the salary (y) at
age 60 ($x = 60$)?



8.1.2 Models as Functions

- Once we have data in an appropriate vector representation, we can construct a **predictive function** (known as a **predictor**).
- Here, a model means a **predictor**.
- A predictor is a function that, when given a particular input example (in our case, a vector of features), produces an output.
- For example,

$$f: \mathbb{R}^D \rightarrow \mathbb{R}$$

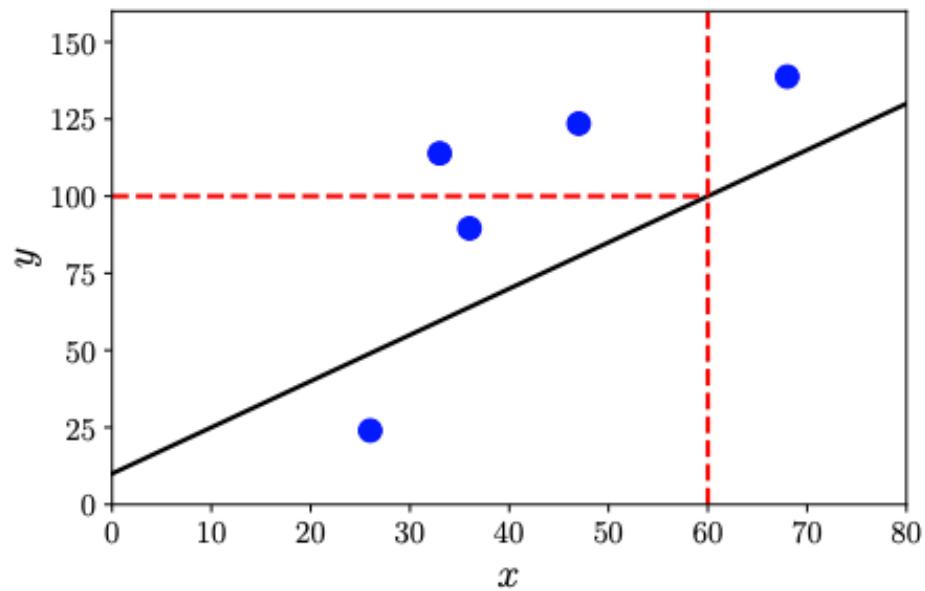
where the input \mathbf{x} is a D -dimensional vector, and the output is a real-valued scalar. That is, the function f is applied to \mathbf{x} , written as $f(\mathbf{x})$ and returns a real number.

8.1.2 Models as Functions

- We mainly consider the special case of linear functions

$$f(\mathbf{x}) = \boldsymbol{\theta}^T \mathbf{x} + \theta_0$$

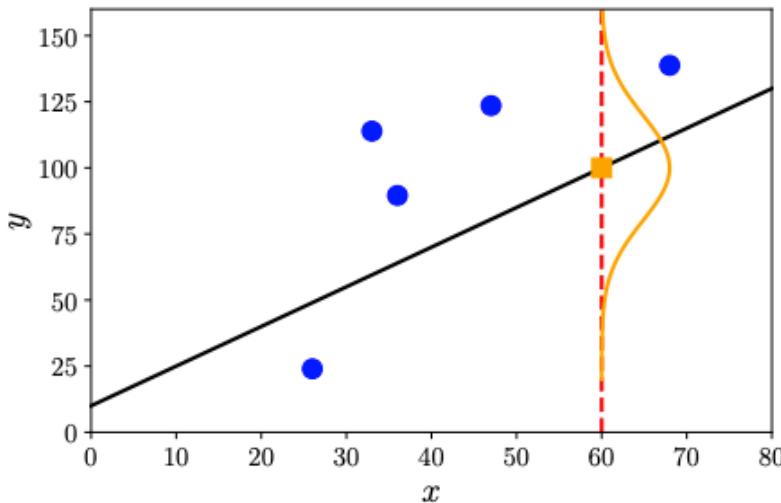
- Example: predicting salary $f(\mathbf{x})$ from age \mathbf{x} .



Black and solid diagonal line
is an example predictor.
 $f(60) = 100$

8.1.3 Models as Probability Distributions

- The observed data is usually a combination of the **true underlying data** and **noise**, i.e., $\tilde{x} = x + n$
- We wish to reveal x from \tilde{x}
- So we would like to have predictors that express some sort of uncertainty, e.g., to quantify the confidence we have about the value of the prediction for a particular test data point.



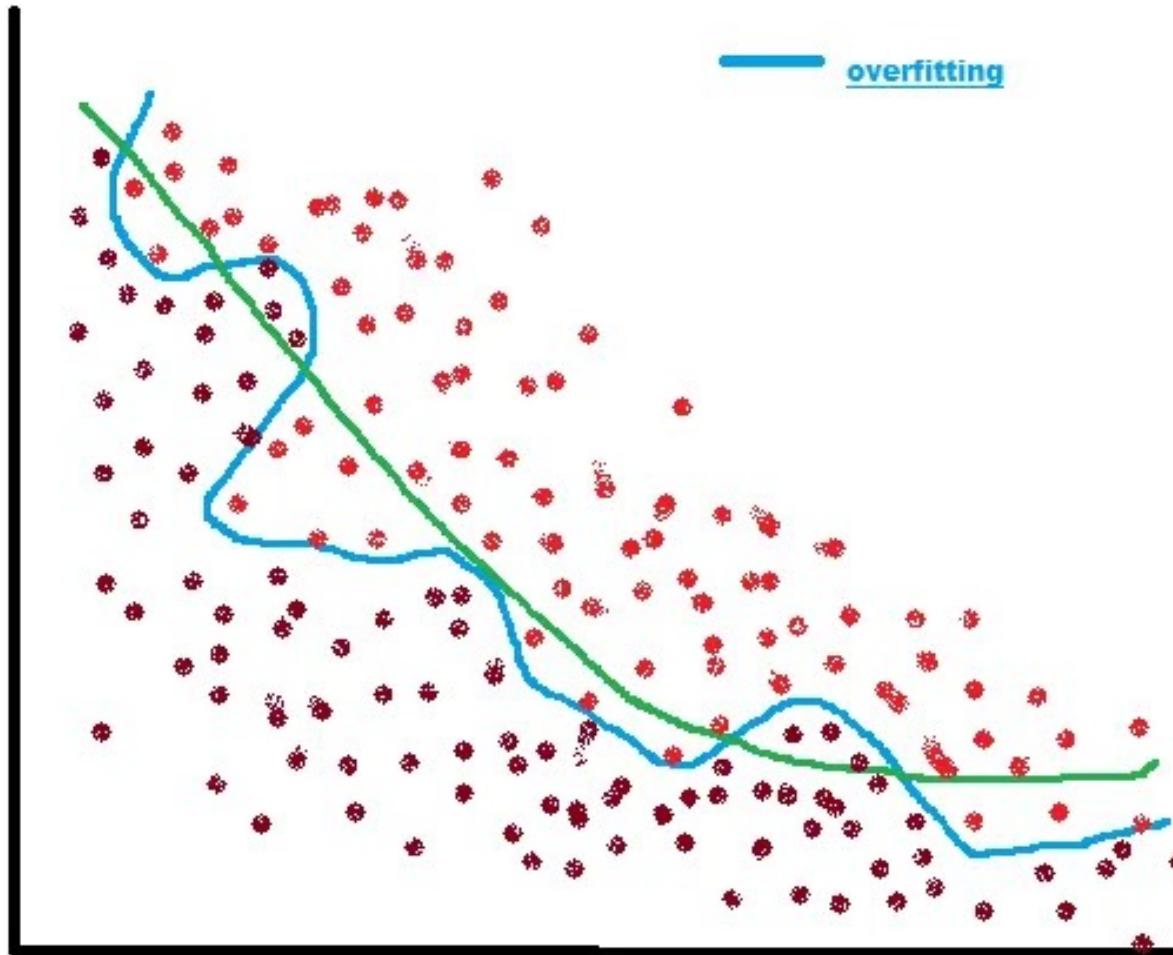
Example function (black solid diagonal line) and its predictive uncertainty at $x = 60$ (drawn as a Gaussian).

- Instead of considering a predictor as a single function, we could consider predictors to be **probabilistic models**.
- We will learn probability in later lectures

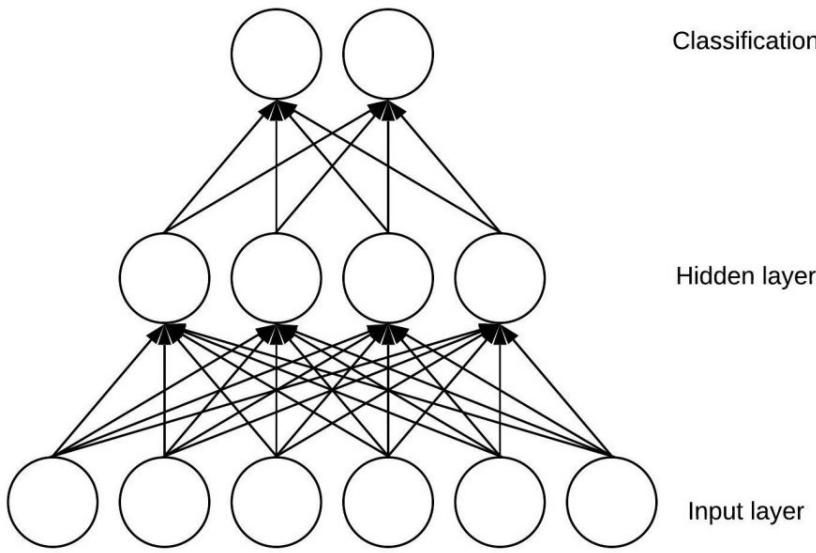
8.1.4 Learning is Finding Parameters

- The goal of learning is to find a **model** and its corresponding **parameters** such that the resulting **predictor** will perform well on unseen **data**.
- 3 algorithmic phases when discussing machine learning algorithms
 - Prediction or inference
 - Training or parameter estimation
 - Hyperparameter tuning or model selection
- Prediction phase: we use a trained predictor on previously unseen test data
- The training or parameter estimation phase: we adjust our predictive model based on training data. We will introduce the **empirical risk minimization** for finding good parameters.
- We use **cross-validation** to assess predictor performance on unseen data.
- We also need to balance between fitting well on training data and finding “simple” explanations of the phenomenon. This trade-off is often achieved using **regularization**.

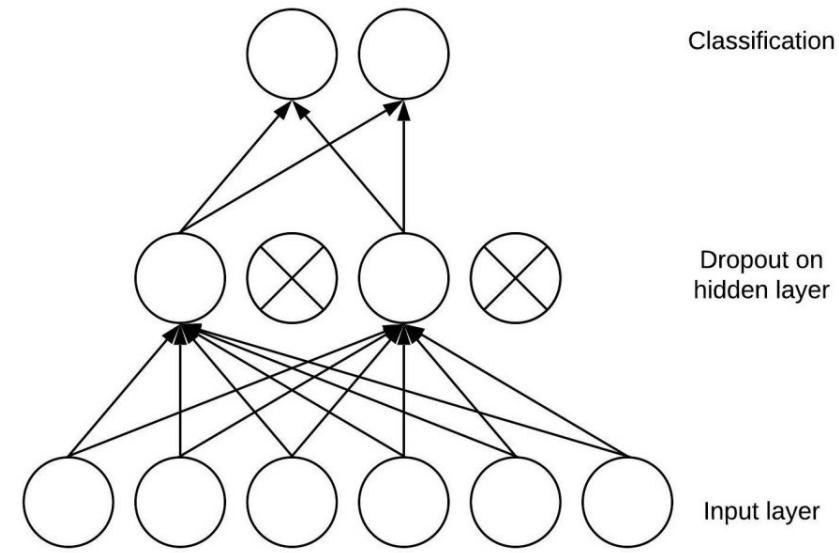
Regularization helps with overfitting.



Regularization in neural networks: dropout.

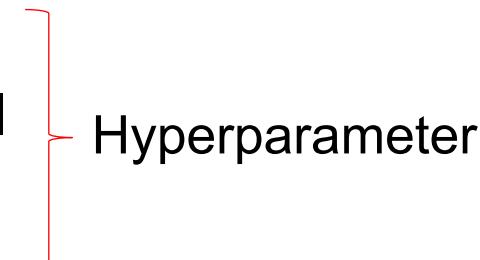


Without Dropout



With Dropout

Credit: Wenkang Wei

- Hyperparameter tuning or model selection
 - We need to make high-level modeling decisions about the structure of the predictor. For example
 - Number of layers to be used in deep learning
 - Number of components in a Gaussian Mixture Model
 - Weight of regularization terms
 - The problem of choosing among different models/hyperparameters is called **model selection**
- Difference between parameters and hyperparameters
- Parameters are to be numerically optimized ($\sim 10^6$ weights in a deep network)
- Hyperparameters need to use search techniques and validation sets (neural architecture search [2]).
- 

Let's take a break

$$m \begin{array}{|c|c|c|c|} \hline & \text{blue} & \text{blue} & \text{blue} \\ \hline \text{blue} & \text{blue} & \text{orange} & \text{blue} \\ \hline & \text{blue} & \text{blue} & \text{blue} \\ \hline \end{array} p = m \begin{array}{|c|c|} \hline \text{orange} & \text{orange} \\ \hline \end{array} \cdot n \begin{array}{|c|c|c|c|} \hline & \text{blue} & \text{blue} & \text{blue} \\ \hline \text{blue} & \text{blue} & \text{orange} & \text{blue} \\ \hline & \text{blue} & \text{blue} & \text{blue} \\ \hline \end{array} p$$

The diagram illustrates matrix multiplication. On the left, a matrix $m \times p$ is shown as a grid of m rows and p columns. One column is highlighted in orange. In the middle, the matrix is multiplied by a scalar n , represented by two orange squares. On the right, the result is a matrix $m \times p$ where each element is multiplied by n . The second column of the result is highlighted in orange, indicating that the scalar multiplication was applied to the second column of the original matrix.

8.2 Empirical Risk Minimization

- What does it mean to **learn**?
- Estimating parameters based on training data.
- Four questions will be answered
- What is the set of functions we allow the predictor to take? –
Hypothesis class of functions
- How do we measure how well the predictor performs on the training data? -- **Loss functions for training**
- How do we construct predictors from only training data that performs well on unseen test data? -- **regularization**
- What is the procedure for searching over the space of models? --
Cross-Validation

8.2.1 Hypothesis Class of Functions

- We are given N examples $\mathbf{x}_n \in \mathbb{R}^D$ and corresponding scalar labels $y_n \in \mathbb{R}$.
- Supervised learning: we have pairs $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$
- We want to estimate a predictor $f(\cdot, \theta) : \mathbb{R}^D \rightarrow \mathbb{R}$, parametrized by θ
- We hope to be able to find a good parameter θ^* such that we fit the data well, that is

$$f(\mathbf{x}_n, \theta^*) \approx y_n \text{ for all } n = 1, \dots, N$$

- We use $\hat{y}_n = f(\mathbf{x}_n, \theta^*)$ to represent the output of the predictor

Example (least-squares regression)

- When the label y_n is real-valued, a popular choice of function class for predictors is **affine functions** (linear functions).

$$f(\mathbf{x}) = \boldsymbol{\theta}^T \mathbf{x} + \theta_0$$

- For more compact representations, we concatenate an additional unit feature $x_n^{(0)} = 1$ to \mathbf{x}_n , i.e.,

$$\mathbf{x}_n = [x_n^{(0)}, x_n^{(1)}, x_n^{(2)}, \dots, x_n^{(D)}]^T = [1, x_n^{(1)}, x_n^{(2)}, \dots, x_n^{(D)}]^T$$

- The parameter vector is $\boldsymbol{\theta} = [\theta_0, \theta_1, \theta_2, \dots, \theta_D]^T$
- We can write the predictor as follows

$$f(\mathbf{x}_n, \boldsymbol{\theta}) = \boldsymbol{\theta}^T \mathbf{x}_n$$

which is equivalent to the affine model

$$f(\mathbf{x}_n, \boldsymbol{\theta}) = \theta_0 + \sum_{d=1}^D \theta_d x_n^{(d)} = \theta_0 x_n^{(0)} + \sum_{d=1}^D \theta_d x_n^{(d)} = \boldsymbol{\theta}^T \mathbf{x}_n$$

Example (least-squares regression)

$$f(\mathbf{x}_n, \boldsymbol{\theta}) = \boldsymbol{\theta}^T \mathbf{x}_n$$

- The predictor takes the vector of features representing a single example \mathbf{x}_n as input and produces a real-valued output,

$$f: \mathbb{R}^{D+1} \rightarrow \mathbb{R}$$

- $f(\mathbf{x}_n, \boldsymbol{\theta}) = \boldsymbol{\theta}^T \mathbf{x}_n$ is a linear predictor
- There are many non-linear predictors, such as the neural networks

8.2.2 Loss Function for Training

- In training, we aim to learn a model that **fits the data well**.
- To define “**fits the data well**”, we specify a **loss function**

$$\ell(y_n, \hat{y}_n)$$

- Input: ground truth label y_n of a training example
the prediction \hat{y}_n of this training example
- Output: a non-negative number, called **loss**. It represents how much error we have made on this particular prediction
- To find good parameters θ^* , we need to minimize the average loss on the set of N training examples
- We usually assume training examples $(x_1, y_1), \dots, (x_N, y_N)$ are **independent and identically distributed (i.i.d.)**.

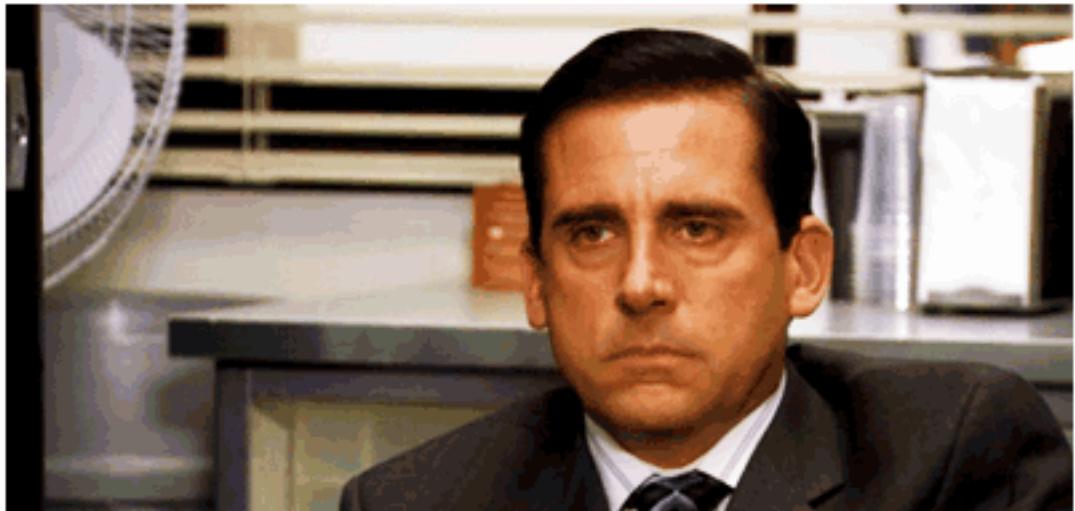
Intuitive example: Learning how to give you the best lecture

5. Which do you prefer?

[More Details](#)

- Slidesss 60
- Blackboard sessions 29
- Hybrid for the win 63
- Other 3

TEACHER MAKES A JOKE
CLASS REACTION:



- Under the i.i.d assumption, the empirical mean is a good estimate of the population mean.
- We can use the empirical mean of the loss on the training data
- Given a **training set** $\{(x_1, y_1), \dots, (x_N, y_N)\}$, we use the notation of an example matrix

$$X := [x_1, \dots, x_N]^T \in \mathbb{R}^{N \times D}$$

and a label vector

$$y = [y_1, \dots, y_N]^T \in \mathbb{R}^N$$

- The average loss is given by

$$R_{emp}(f, X, y) = \frac{1}{N} \sum_{n=1}^N \ell(y_n, \hat{y}_n)$$

where $\hat{y}_n = f(x_n, \theta)$. The above equation is called the **empirical risk**. The learning strategy is called **empirical risk minimization**.

Example - Least-Squares Loss

- We use the squared loss function

$$\ell(y_n, \hat{y}_n) = (y_n - \hat{y}_n)^2$$

- We aim to minimize the empirical risk, which is the average of the losses over the training data.

$$\min_{\theta \in \mathbb{R}^D} \frac{1}{N} \sum_{n=1}^N \ell(y_n, \hat{y}_n) = \min_{\theta \in \mathbb{R}^D} \frac{1}{N} \sum_{n=1}^N (y_n - f(\mathbf{x}_n, \theta))^2$$

Using the linear predictor $f(\mathbf{x}_n, \theta) = \theta^\top \mathbf{x}_n$, we obtain the optimization problem

$$\min_{\theta \in \mathbb{R}^D} \frac{1}{N} \sum_{n=1}^N (y_n - f(\mathbf{x}_n, \theta))^2$$

- This equation can be equivalently expressed in matrix form

$$\min_{\theta \in \mathbb{R}^D} \frac{1}{N} \|y - X\theta\|^2$$

- This is known as the **least-squares problem**. There exists a closed-form analytic solution for this by solving the normal equations. We will discuss it in later lectures

- We actually want to find a predictor f that minimizes the **expected risk** (or the population risk)

$$R_{\text{true}}(f) = \mathbb{E}_{x,y}[\ell(y, f(x))]$$

where y is the ground truth label and $f(x)$ is the prediction based on the example x .

- $R_{\text{true}}(f)$ is the true risk, if we can access an infinite amount of data
- The expectation \mathbb{E} is over the infinite set of all possible data and labels.

- Machine learning applications have different types of performance measure.
 - For classification: accuracy, AUC, F1 score, etc.
 - For detection: mean average precision, mIoU, etc.
 - For image denoise/super resolution: SSIM, PSNR, etc.
- In principle, the loss function should correspond to the measure.
- However, there are often mismatches between loss functions and the measures – due to implementation/optimization considerations.

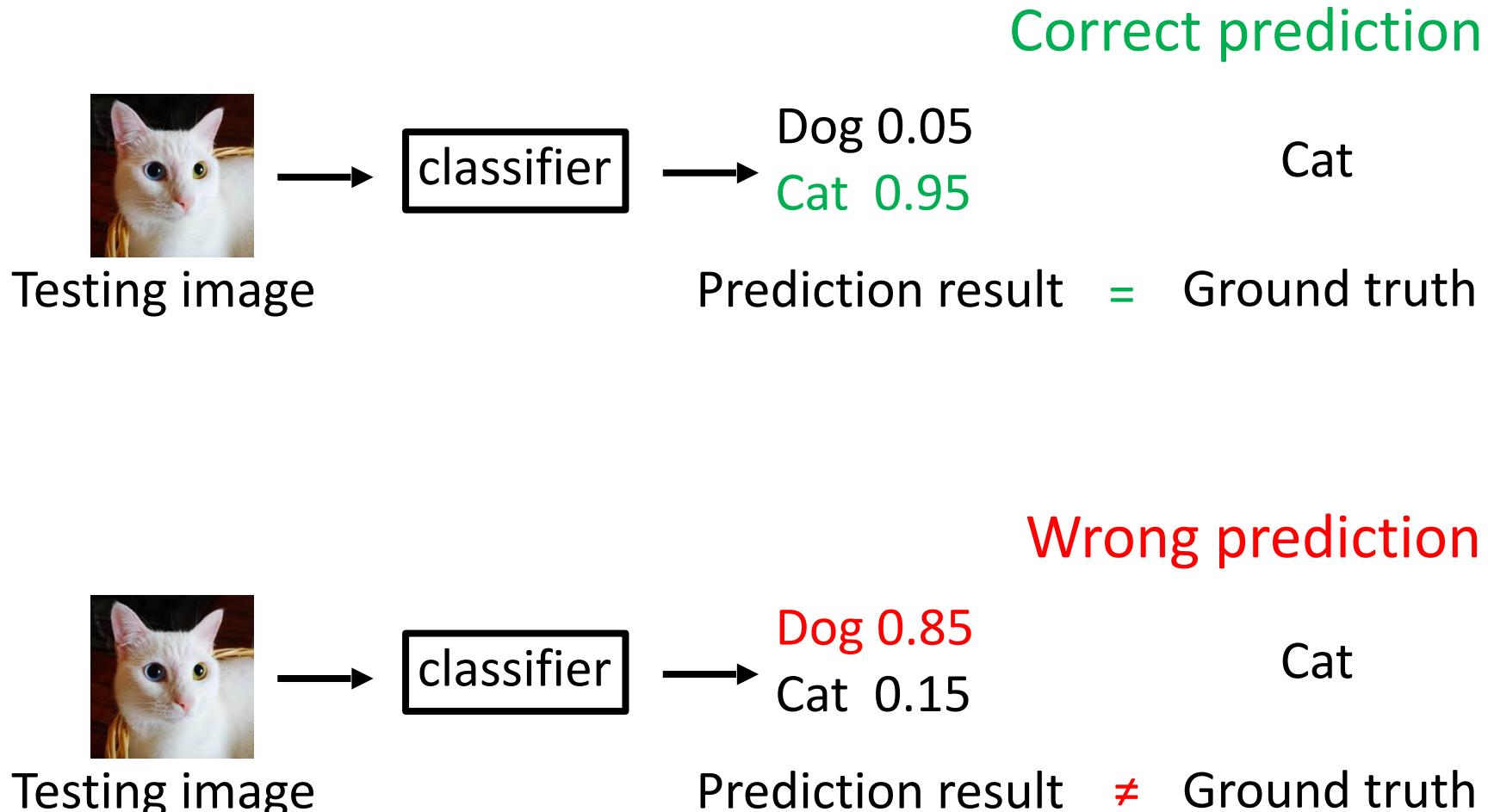
We start with training a classifier



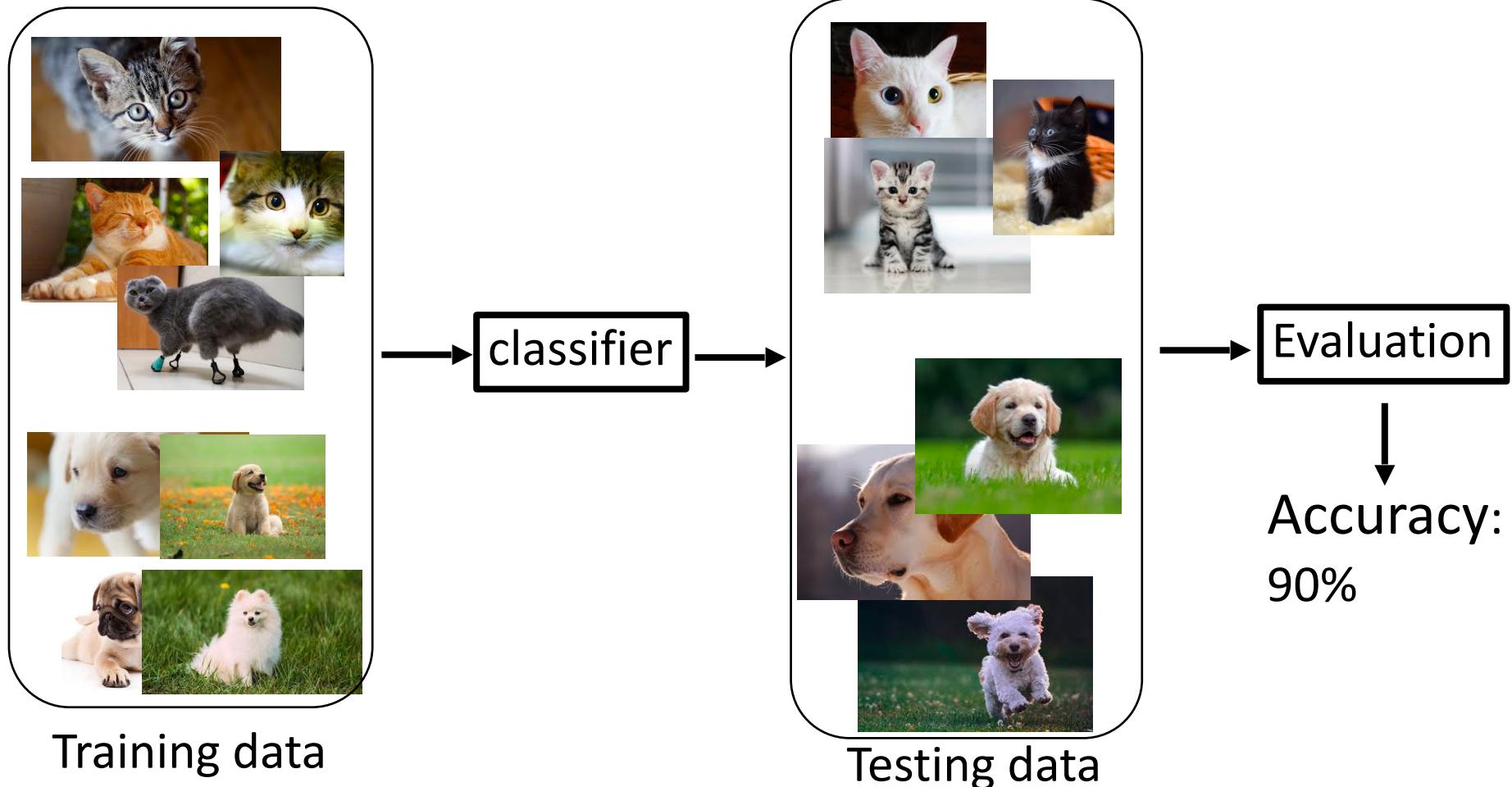
→ classifier

Training data

We do a bit testing....



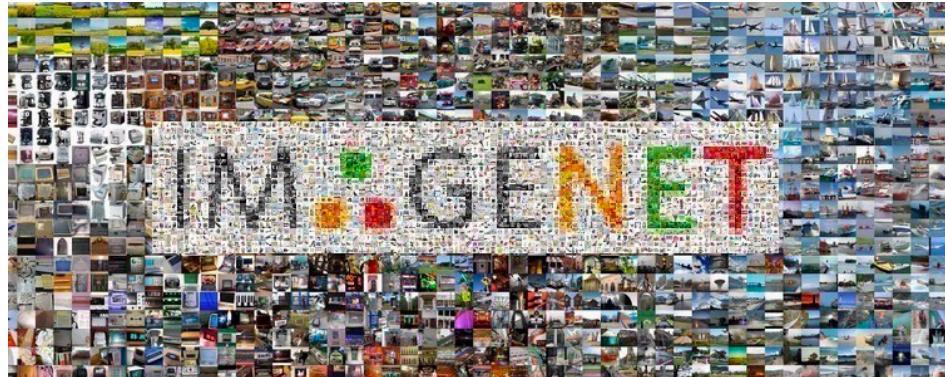
We now evaluate a model



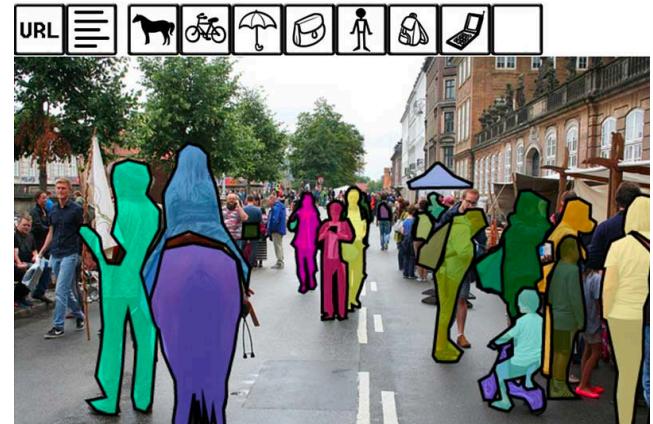
Ground truths provided

Is this way of evaluation feasible?

- Yes



ImageNet



MSCOCO

Ground truths provided



LFW

Is this way of evaluation feasible?

- No....

We can't calculate a classifier accuracy!!

Suppose we deploy our cat-dog classifier to a swimming pool



Ground truths not provided

Research Bit: finding stars that weren't born in our Galaxy

