# Deep Learning for Data Science DS 542

https://dl4ds.github.io/sp2026/

Supervised Learning

# Supervised learning

- Examples
- Terminology
- Notation
  - Model
  - Loss function
  - Training
  - Testing
- 1D Linear regression example
  - Model
  - Loss function
  - Training
  - Testing

# Supervised learning

- Define a mapping from input to output
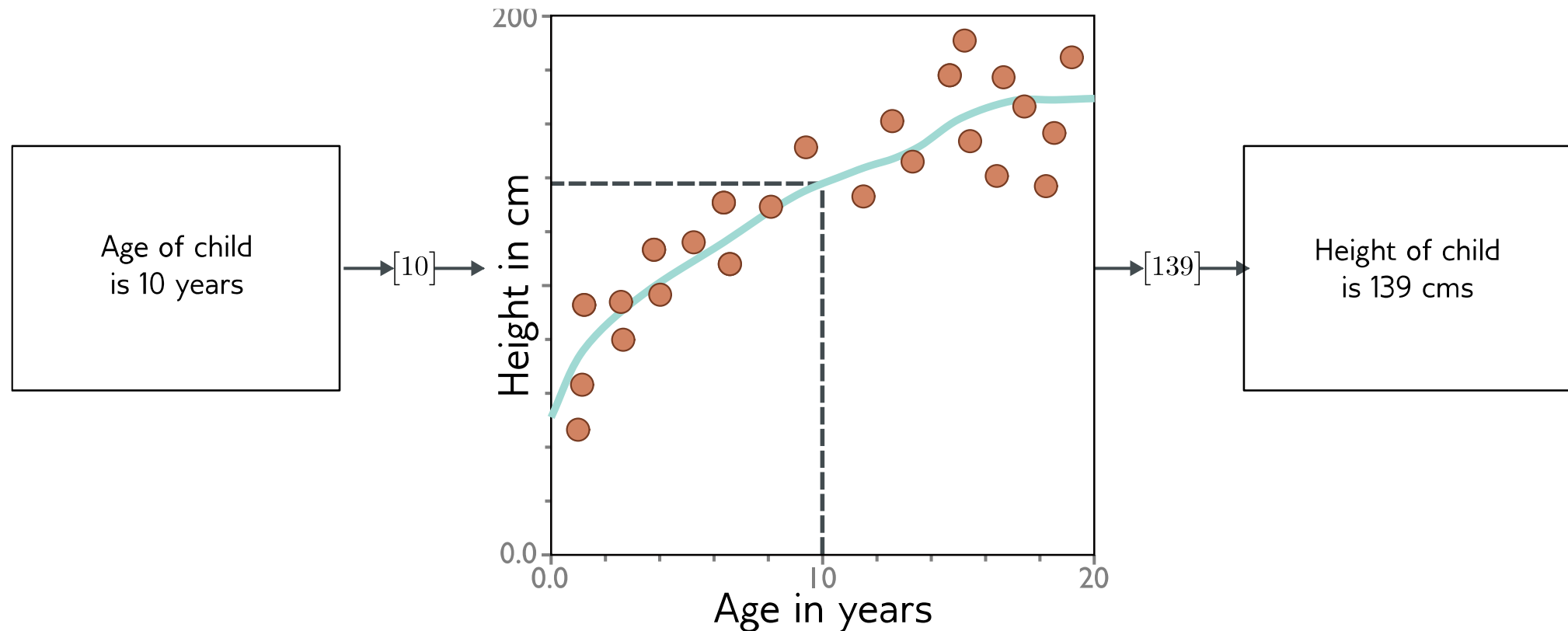
- Learn this mapping from paired input/output data examples



Parameter Updates

Input Data

Parameterized Model

Output Predictions

Loss Calc & Parameter Updates
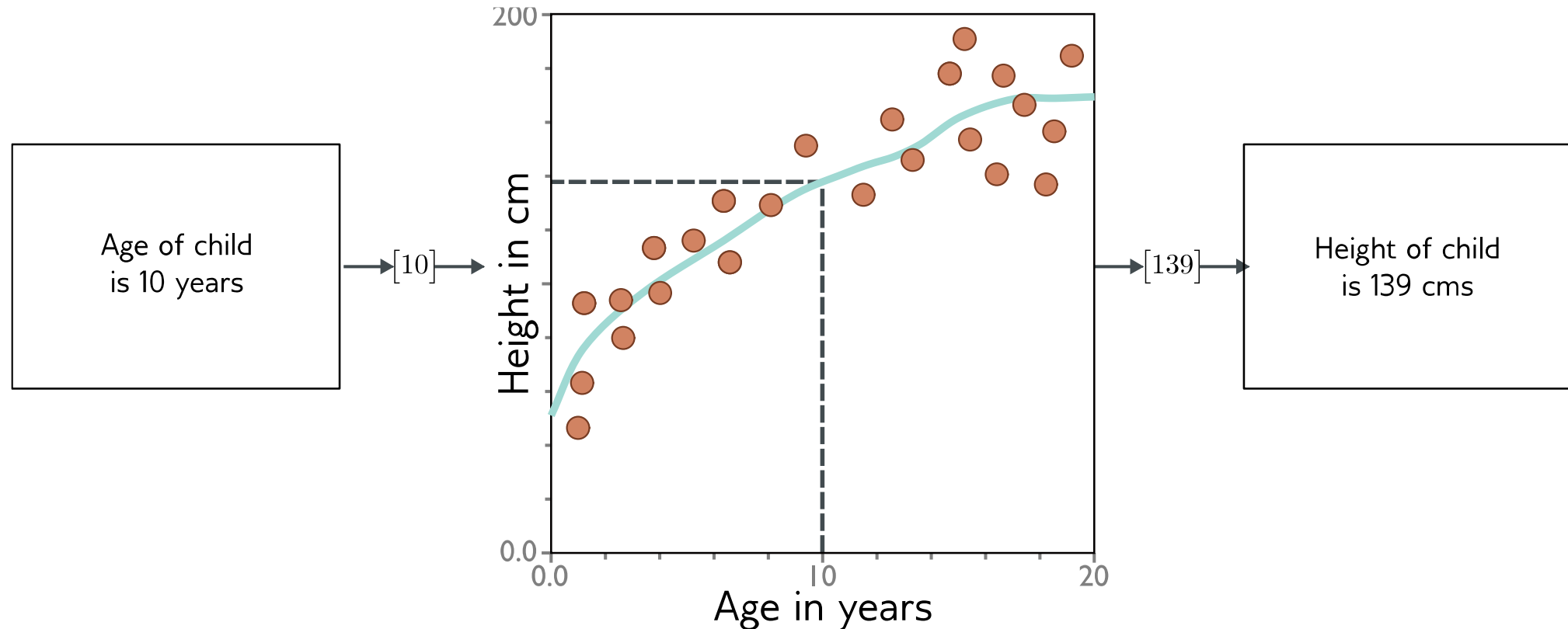
Ground Truth / Labels / Targets

# What is a supervised learning model?



- An equation relating input (age) to output (height)
- Search through family of possible equations to find one that fits training data well

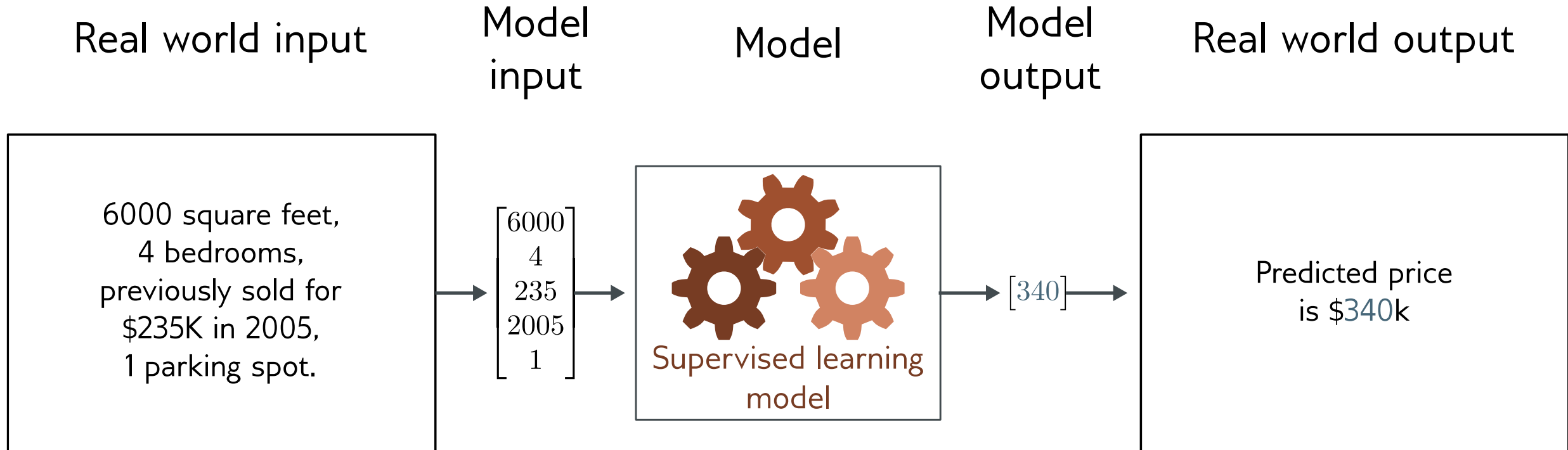# What is a supervised learning model?



- Deep neural networks are just a very flexible family of equations
- Fitting deep neural networks = "Deep Learning"

# Prediction Types

- Regression
  - Prediction a continuous valued output

- Classification
  - Assigning input to one of a finite number of classes or categories
  - Two classes are a special case

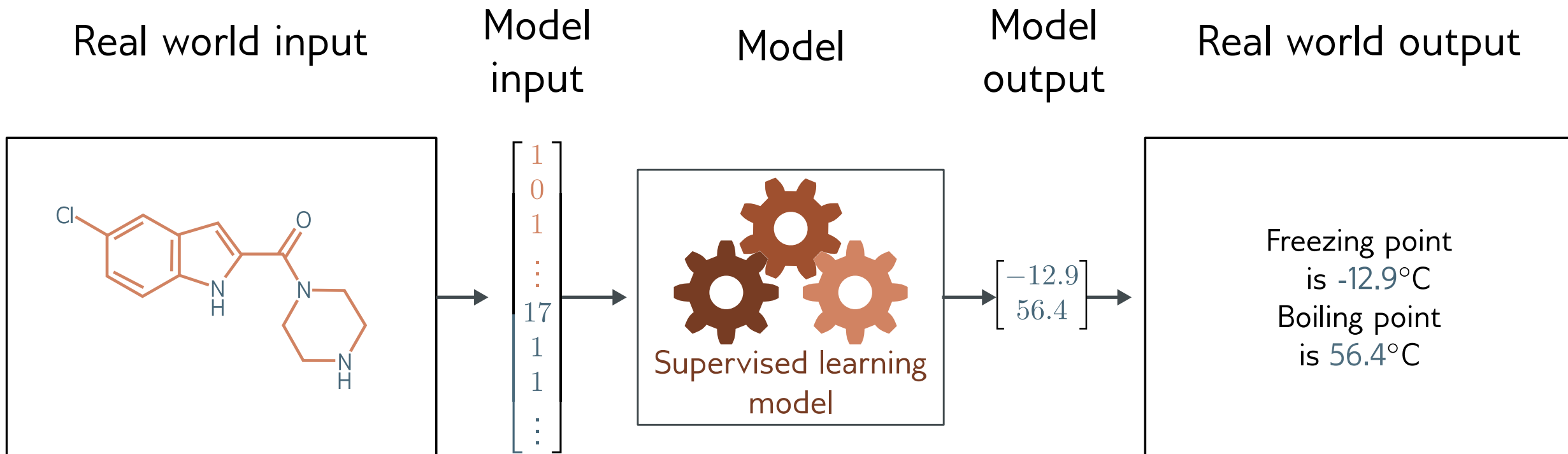  Can be univariate (one output) or multivariate ( more than one output)

# Regression

Real world input

Model input

Model

Model output

Real world output

6000 square feet,
4 bedrooms,
previously sold for
$235K in 2005,
1 parking spot.

$$\begin{bmatrix} 6000 \\ 4 \\ 235 \\ 2005 \\ 1 \end{bmatrix}$$

Supervised learning model

$[340]$

Predicted price
is $340k

- Univariate regression problem (one output, real value)
- Fully connected network

# Graph regression

Real world input

Model
input

Model

Model
output

Real world output

$$\begin{bmatrix} 1 \\ 0 \\ 1 \\ \vdots \\ 17 \\ 1 \\ 1 \\ 1 \\ \vdots \end{bmatrix}$$

Supervised learning
model

$$\begin{bmatrix} -12.9 \\ 56.4 \end{bmatrix}$$

Freezing point
is -12.9°C
Boiling point
is 56.4°C

- Multivariate regression problem (>1 output, real value)
- Graph neural network

# Text classification

Real world input | Model input | Model | Model output | Real world output

"The steak was terrible, the salad was rotten, and the soup tasted like socks"

$$\begin{bmatrix} 8672 \\ 8194 \\ 9804 \\ 8634 \\ 8672 \\ \vdots \end{bmatrix}$$

Supervised learning model

$$\begin{bmatrix} 0.02 \\ 0.98 \end{bmatrix}$$

Positive
Negative

- Binary classification problem (two discrete classes)
- Transformer network

# Music genre classification

Real world input Model input Model Model output Real world output



$$\begin{bmatrix} 125 \\ 12054 \\ 1253 \\ 6178 \\ 24 \\ 4447 \\ \vdots \end{bmatrix}$$

Supervised learning model

$$\begin{bmatrix} 0.03 \\ 0.52 \\ 0.18 \\ 0.07 \\ 0.12 \\ 0.08 \\ \vdots \\ 0.01 \end{bmatrix}$$
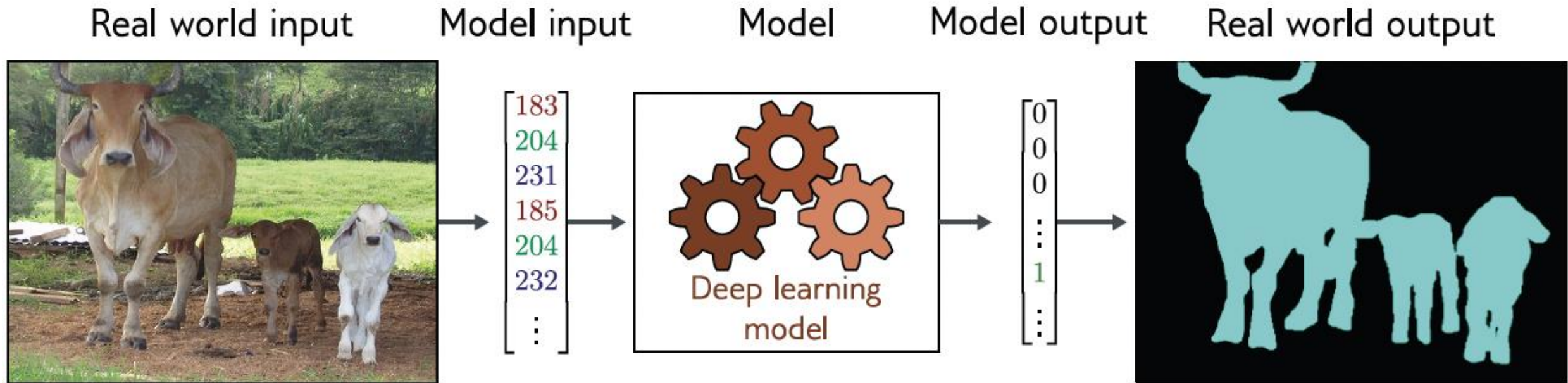
Classical
Electronica
Hip Hop
Jazz
Pop
Metal
Punk

- Multiclass classification problem (discrete classes, >2 possible values)
- Recurrent neural network (RNN)

# Image classification

Real world input
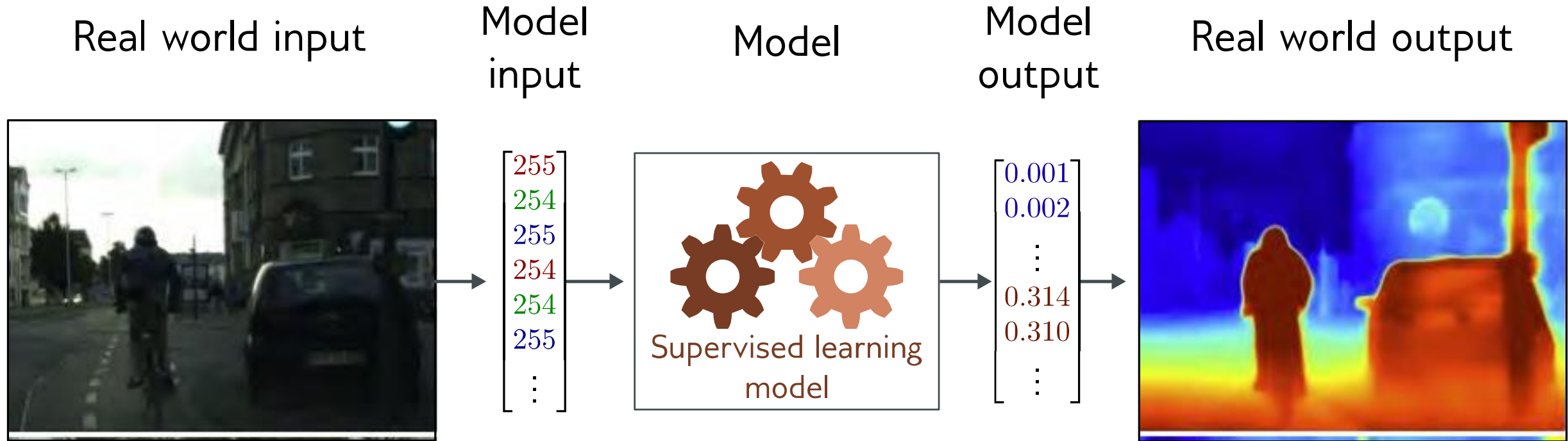
Model input

Model

Model output

Real world output

$$\begin{bmatrix} 124 \\ 140 \\ 156 \\ 128 \\ 142 \\ 157 \\ \vdots \end{bmatrix}$$

Supervised learning model

$$\begin{bmatrix} 0.00 \\ 0.00 \\ 0.01 \\ 0.89 \\ 0.05 \\ 0.00 \\ \vdots \\ 0\ 01 \end{bmatrix}$$

Aardvark
Apple
Bee
Bicycle
Bridge
Clown
⋮

- Multiclass classification problem (discrete classes, >2 possible classes)
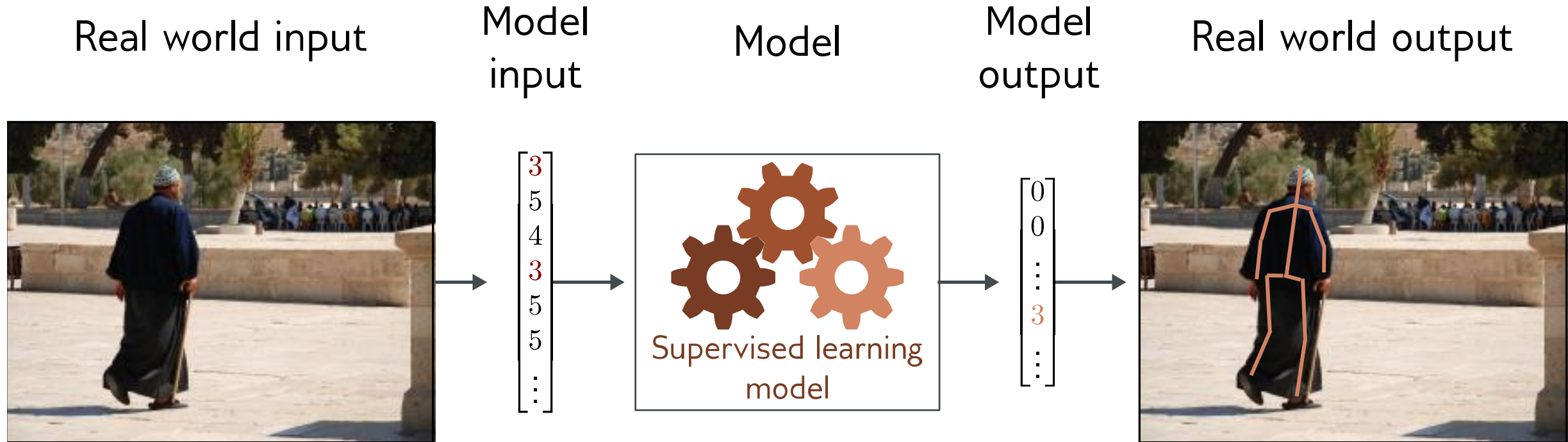- Convolutional network

# Image segmentation



Real world input  Model input  Model  Model output  Real world output

$$\begin{bmatrix} 183 \\ 204 \\ 231 \\ 185 \\ 204 \\ 232 \\ \vdots \end{bmatrix} \qquad \text{Deep learning model} \qquad \begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \end{bmatrix}$$

- Multivariate binary classification problem (many outputs, two discrete classes)
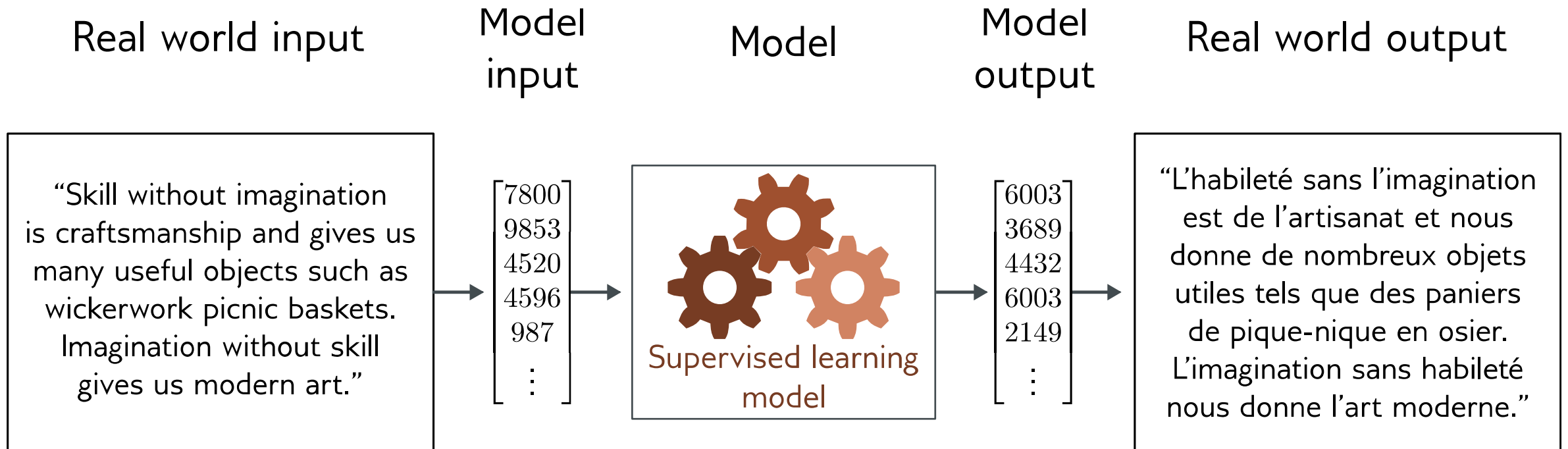- Convolutional encoder-decoder network

# Depth estimation

Real world input    Model input    Model    Model output    Real world output



$$\begin{bmatrix} 255 \\ 254 \\ 255 \\ 254 \\ 254 \\ 255 \\ \vdots \end{bmatrix}$$

Supervised learning model

$$\begin{bmatrix} 0.001 \\ 0.002 \\ \vdots \\ 0.314 \\ 0.310 \\ \vdots \end{bmatrix}$$

- Multivariate regression problem (many outputs, continuous)
- Convolutional encoder-decoder network

# Pose estimation

Real world input

Model input

Model

Model output

Real world output



Supervised learning model

- Multivariate regression problem (many outputs, continuous)
- Convolutional encoder-decoder network

# Translation

Real world input

Model

Model
output

Real world output

"Skill without imagination is craftsmanship and gives us many useful objects such as wickerwork picnic baskets. Imagination without skill gives us modern art."

$$\begin{bmatrix} 7800 \\ 9853 \\ 4520 \\ 4596 \\ 987 \\ \vdots \end{bmatrix}$$



Supervised learning model

$$\begin{bmatrix} 6003 \\ 3689 \\ 4432 \\ 6003 \\ 2149 \\ \vdots \end{bmatrix}$$

"L'habileté sans l'imagination est de l'artisanat et nous donne de nombreux objets utiles tels que des paniers de pique-nique en osier. L'imagination sans habileté nous donne l'art moderne."

- Encoder-Decoder Transformer Networks

UDL    16

# Image captioning

Real world input

Model input

Model

Model output

Real world output



$\begin{bmatrix} 183 \\ 204 \\ 231 \\ 185 \\ 204 \\ 232 \\ \vdots \end{bmatrix}$

Supervised learning model

$\begin{bmatrix} 1 \\ 5593 \\ 7532 \\ 7924 \\ 1 \\ \vdots \end{bmatrix}$

"A Kazakh man on a horse holding a bird of prey"

- E.g. CNN-RNN, LSTM, Transformers

# Image generation from text

Real world input | Model input | Model | Model output | Real world output



"Teddy bears mixing sparkling chemicals as mad scientists, in a steampunk style."

$$\begin{bmatrix} 8300 \\ 532 \\ 7676 \\ 7898 \\ 883 \\ \vdots \end{bmatrix}$$

Supervised learning model

$$\begin{bmatrix} 2 \\ 2 \\ 2 \\ 2 \\ 2 \\ 1 \\ \vdots \end{bmatrix}$$

# Supervised Learning Classification and Regression Applications

# Regression

| Real world input | Model input | Model | Model output | Real world output |
|---|---|---|---|---|

6000 square feet,
4 bedrooms,
previously sold for
$235K in 2005,
1 parking spot.

$$\begin{bmatrix} 6000 \\ 4 \\ 235 \\ 2005 \\ 1 \end{bmatrix}$$

Supervised learning
model

$[340]$

Predicted price
is $340k

- Univariate regression problem (one output, real value)

# Any Questions?

# Supervised learning

- Examples
- Terminology
- Notation
  - Model
  - Loss function
  - Training
  - Testing
- 1D Linear regression example
  - Model
  - Loss function
  - Training
  - Testing

# Supervised learning terminology

- Supervised learning model = mapping from one or more inputs to one or more outputs

- Model is a family of equations → "inductive bias"

- Computing the outputs from the inputs → inference

- Model also includes parameters

- Parameters affect outcome of equation

- Training a model = finding parameters that predict outputs "well" from inputs for training and evaluation datasets of input/output pairs

# Supervised learning

- Examples
- Terminology
- Notation
  - Model
  - Loss function
  - Training
  - Testing
- 1D Linear regression example
  - Model
  - Loss function
  - Training
  - Testing

# Notation:

- Input:

$$\mathbf{x}$$

Variables always Roman letters

Normal lower case = scalar
Bold lower case = vector
Capital Bold = matrix

- Output:

$$\mathbf{y}$$

- Model:

$$\mathbf{y} = \mathbf{f}[\mathbf{x}]$$

Functions always square brackets

Normal lower case = returns scalar
Bold lower case = returns vector
Capital Bold = returns matrix

# Notation example:

- Input:

$$\mathbf{x} = \begin{bmatrix} \text{age} \\ \text{mileage} \end{bmatrix}$$

<span style="color:#C87A5A">Vector:<br>Structured</span> or<br><span style="color:#C87A5A">tabular</span> data

- Output:

$$y = \begin{bmatrix} \text{price} \end{bmatrix}$$

<span style="color:#C87A5A">Scalar output</span>

- Model:

$$y = \text{f}[\mathbf{x}]$$

<span style="color:#C87A5A">Scalar output<br>function<br>(with vector input)</span>

# Model

- Parameters:

$$\phi$$

- Model :

$$\mathbf{y} = \mathbf{f}[\mathbf{x}, \phi]$$

# Data Set and Loss function

- Training dataset of $I$ pairs of input/output examples:

$$\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^{I}$$

# Data Set and Loss function

- Training dataset of $I$ pairs of input/output examples:

$$\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^{I}$$

- Loss function or cost function measures how bad model is:

$$L\left[\boldsymbol{\phi}, \underbrace{\mathrm{f}[\mathbf{x}, \boldsymbol{\phi}]}_{\text{model}}, \underbrace{\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^{I}}_{\text{train data}}\right]$$

# Data Set and Loss function

- Training dataset of $I$ pairs of input/output examples:

$$\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^{I}$$

- Loss function or cost function measures how bad model is:

$$L\left[\boldsymbol{\phi}, \mathrm{f}[\mathbf{x}, \boldsymbol{\phi}], \{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^{I}\right]$$

$$\underbrace{\phantom{\mathrm{f}[\mathbf{x}, \boldsymbol{\phi}]}}_{\text{model}} \underbrace{\phantom{\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^{I}}}_{\text{train data}}$$

or for short:

$$L\left[\boldsymbol{\phi}\right]$$

Returns a scalar that is smaller when model maps inputs to outputs better

# Training

- Loss function:

$$L\left[\phi\right]$$ ← Returns a scalar that is smaller when model maps inputs to outputs better

- Find the parameters that minimize the loss:

$$\hat{\phi} = \underset{\phi}{\operatorname{argmin}}\left[\mathrm{L}\left[\phi\right]\right]$$

# Any Questions?

# Supervised learning

- Examples
- Terminology
- Notation
  - Model
  - Loss function
  - Training
  - Testing
- 1D Linear regression example
  - Model
  - Loss function
  - Training
  - Testing

# Example: 1D Linear regression model

- Model:

$$y = \mathrm{f}[x, \boldsymbol{\phi}]$$

$$= \phi_0 + \phi_1 x$$

- Parameters

$$\boldsymbol{\phi} = \begin{bmatrix} \phi_0 \\ \phi_1 \end{bmatrix}$$

$\longleftarrow$ y-offset

$\longleftarrow$ slope

# Example: 1D Linear regression model

- Model:

$$y = \mathrm{f}[x, \boldsymbol{\phi}]$$

$$= \phi_0 + \phi_1 x$$

- Parameters

$$\boldsymbol{\phi} = \begin{bmatrix} \phi_0 \\ \phi_1 \end{bmatrix} \quad \longleftarrow \text{ y-offset}$$
$$\longleftarrow \text{ slope}$$

# Example: 1D Linear regression model

- Model:

$$y = \mathrm{f}[x, \boldsymbol{\phi}]$$

$$= \phi_0 + \phi_1 x$$

- Parameters

$$\boldsymbol{\phi} = \begin{bmatrix} \phi_0 \\ \phi_1 \end{bmatrix} \quad \longleftarrow \quad \text{y-offset}$$
$$\longleftarrow \quad \text{slope}$$

# Example: 1D Linear regression model

- Model:

$$y = \mathrm{f}[x, \boldsymbol{\phi}]$$

$$= \phi_0 + \phi_1 x$$

- Parameters

$$\boldsymbol{\phi} = \begin{bmatrix} \phi_0 \\ \phi_1 \end{bmatrix}$$

$\longleftarrow$ y-offset

$\longleftarrow$ slope

Interactive Figure 2.1

# Example: 1D Linear regression training data

# Example: 1D Linear regression training data



Loss function:

$$L[\phi] = \sum_{i=1}^{I} (\mathrm{f}[x_i, \phi] - y_i)^2$$

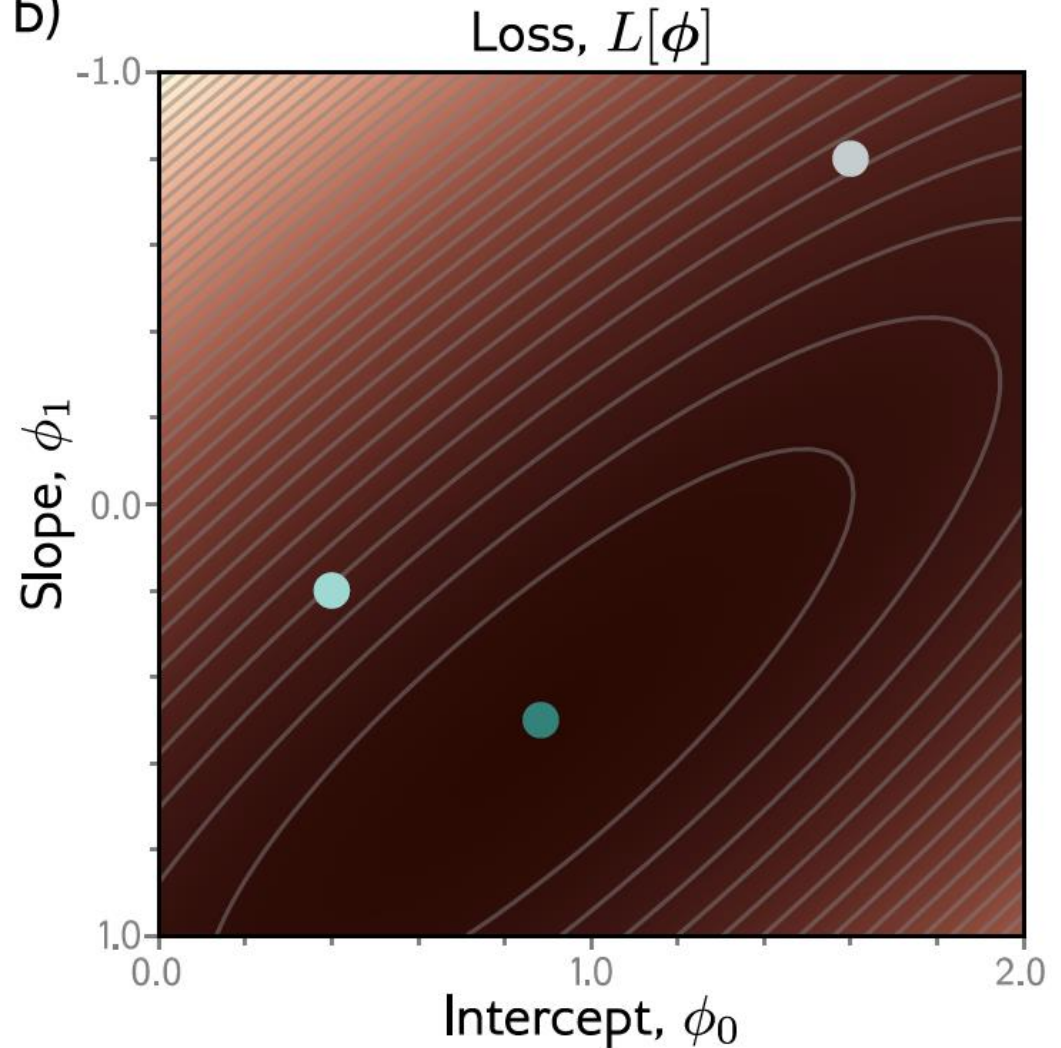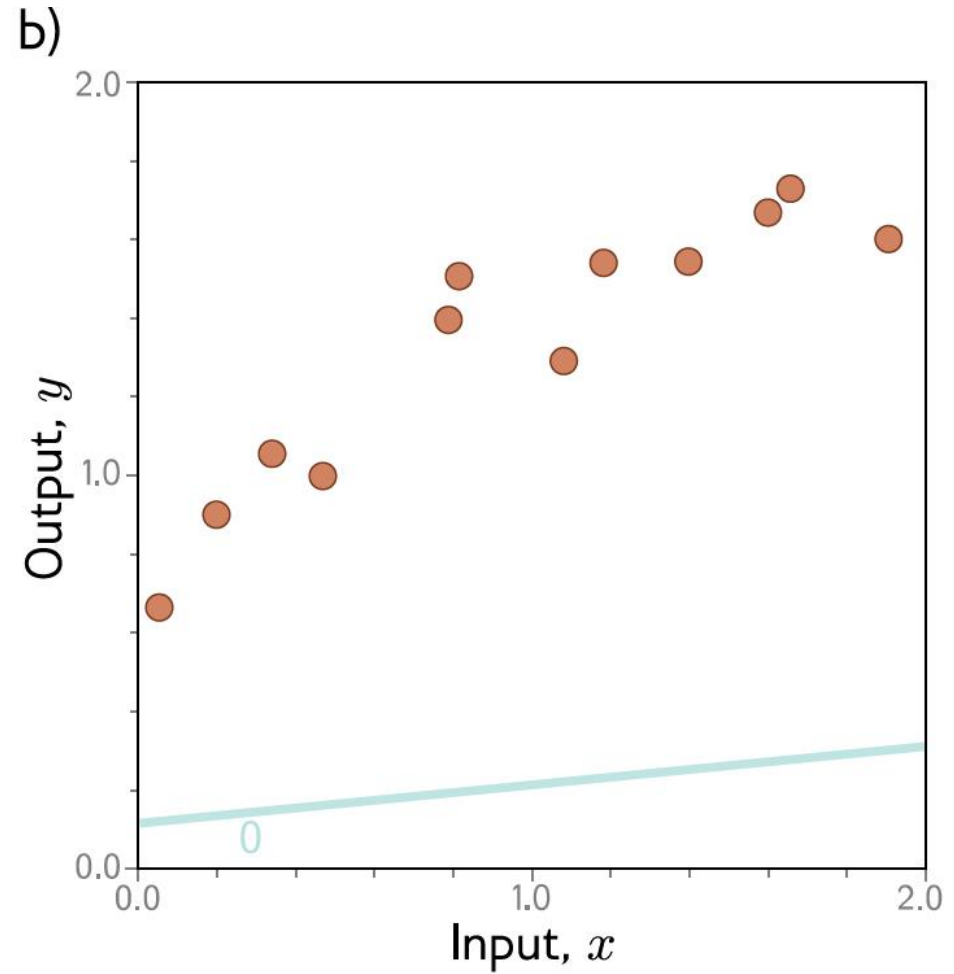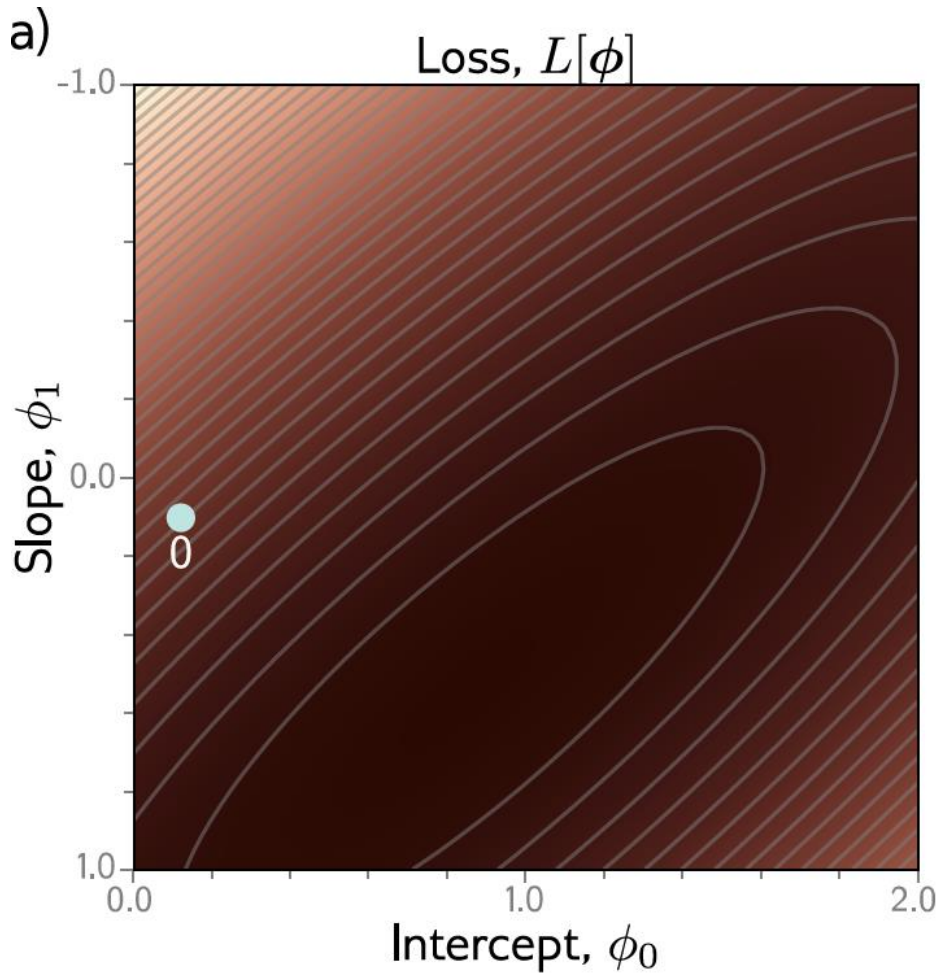$$= \sum_{i=1}^{I} (\phi_0 + \phi_1 x_i - y_i)^2$$

"Least squares loss function"
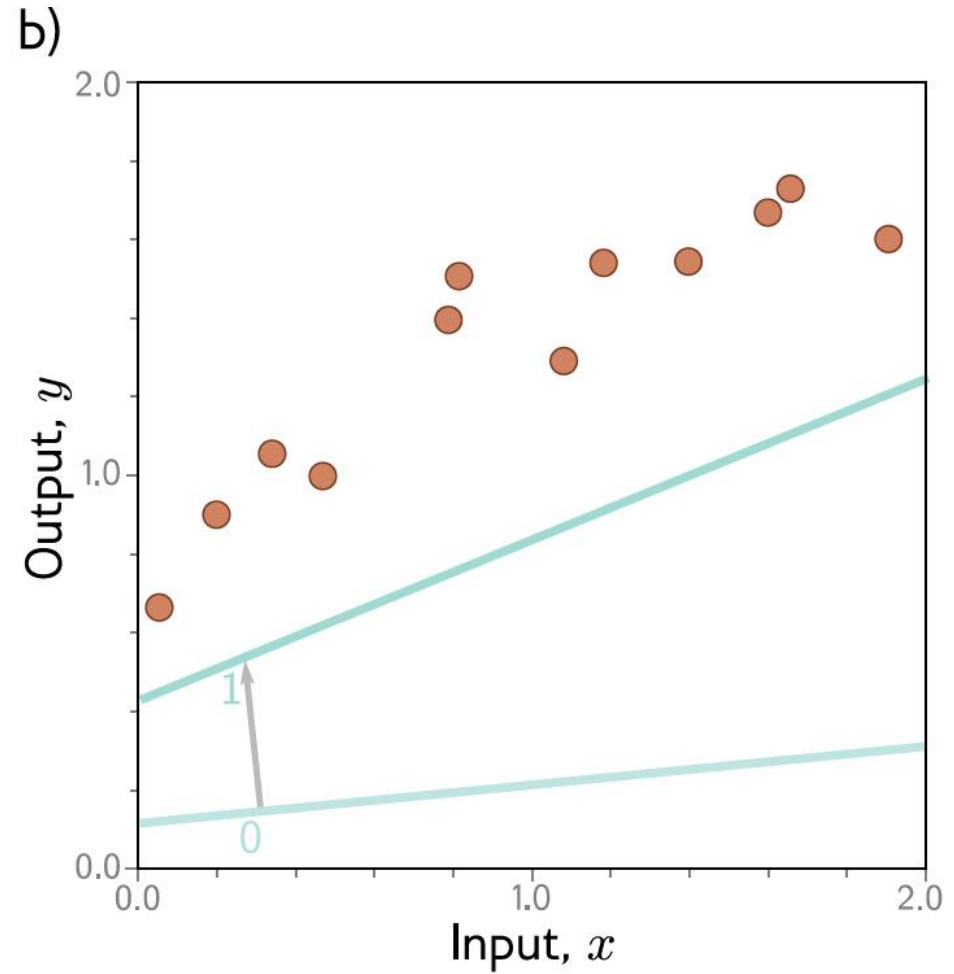
# Example: 1D Linear regression loss function

Loss, $L = 7.11$

Output, $y$

$\phi_0 = 0.4, \phi_1 = 0.2$

Input, $x$

Loss function:

$$L[\boldsymbol{\phi}] = \sum_{i=1}^{I} (\mathrm{f}[x_i, \boldsymbol{\phi}] - y_i)^2$$

$$= \sum_{i=1}^{I} (\phi_0 + \phi_1 x_i - y_i)^2$$

"Least squares loss function"

# Example: 1D Linear regression loss function



Loss, $L = 10.22$

$\phi_0 = 1.60, \phi_1 = -0.8$

Output, $y$

Input, $x$

Loss function:

$$L[\boldsymbol{\phi}] = \sum_{i=1}^{I} (\mathrm{f}[x_i, \boldsymbol{\phi}] - y_i)^2$$

$$= \sum_{i=1}^{I} (\phi_0 + \phi_1 x_i - y_i)^2$$

"Least squares loss function"

# Example: 1D Linear regression loss function



Loss function:

$$L[\boldsymbol{\phi}] = \sum_{i=1}^{I} (\mathrm{f}[x_i, \boldsymbol{\phi}] - y_i)^2$$
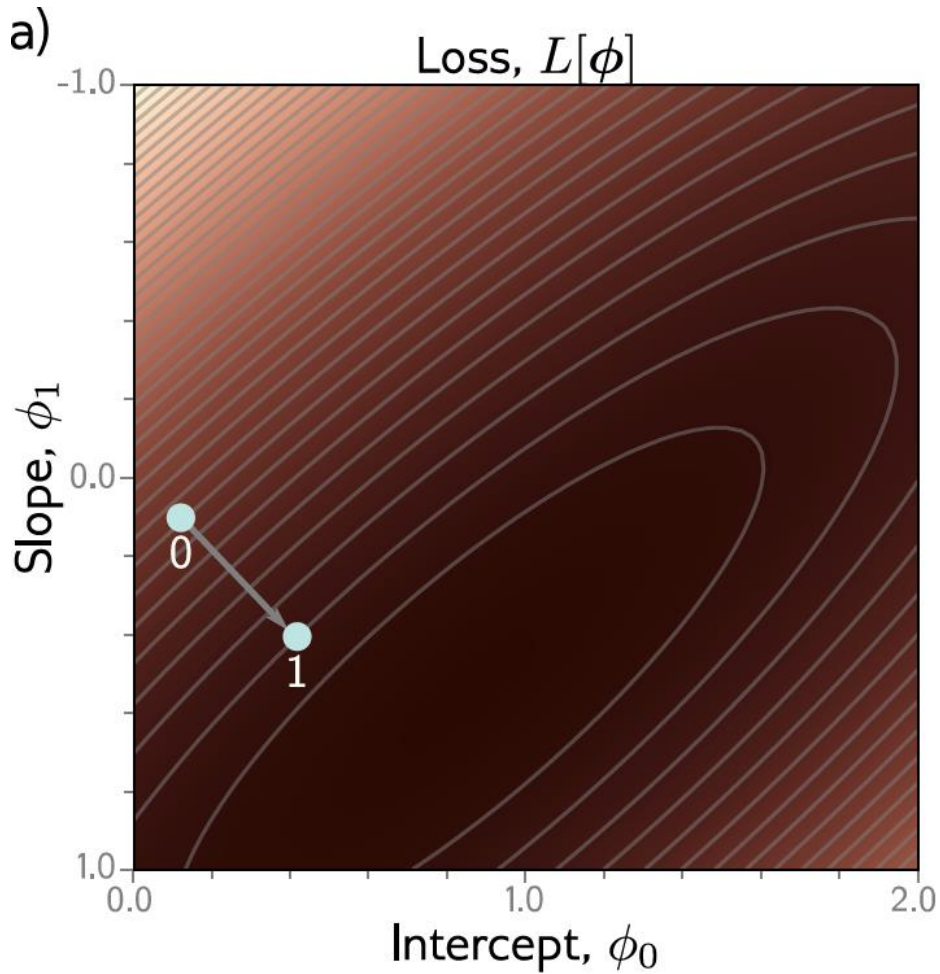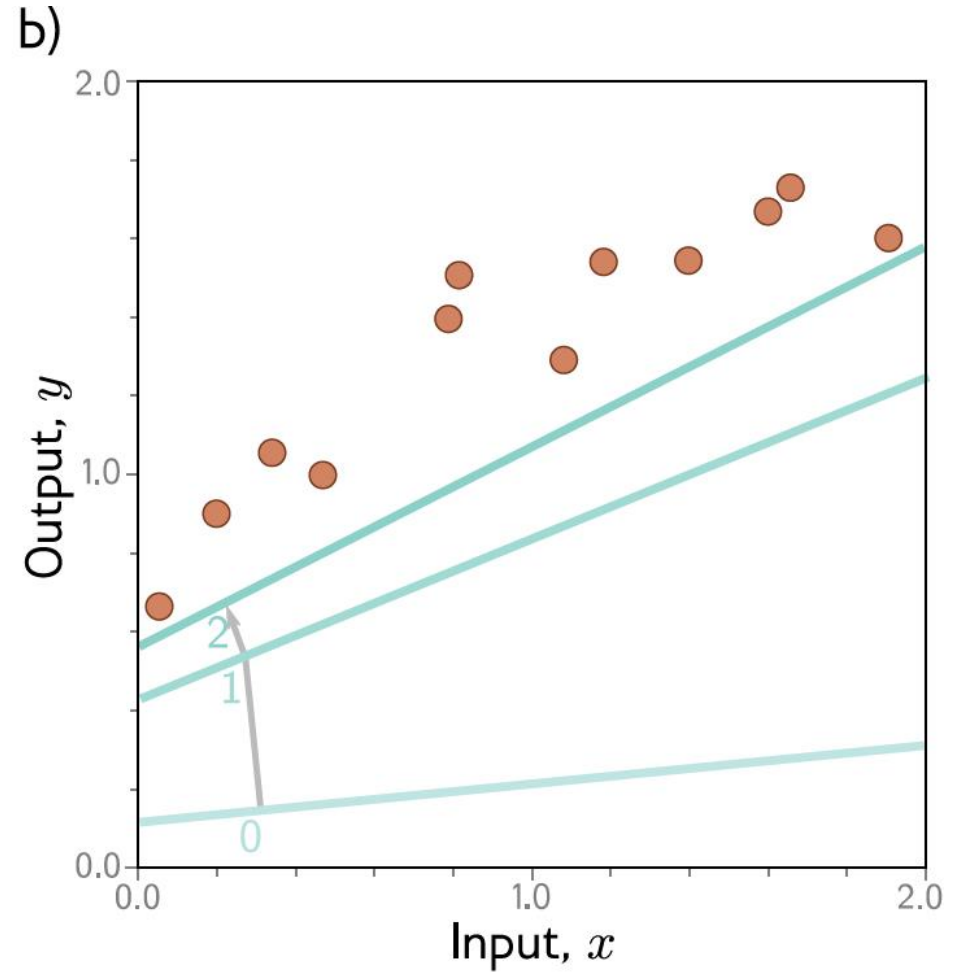
$$= \sum_{i=1}^{I} (\phi_0 + \phi_1 x_i - y_i)^2$$

"Least squares loss function"
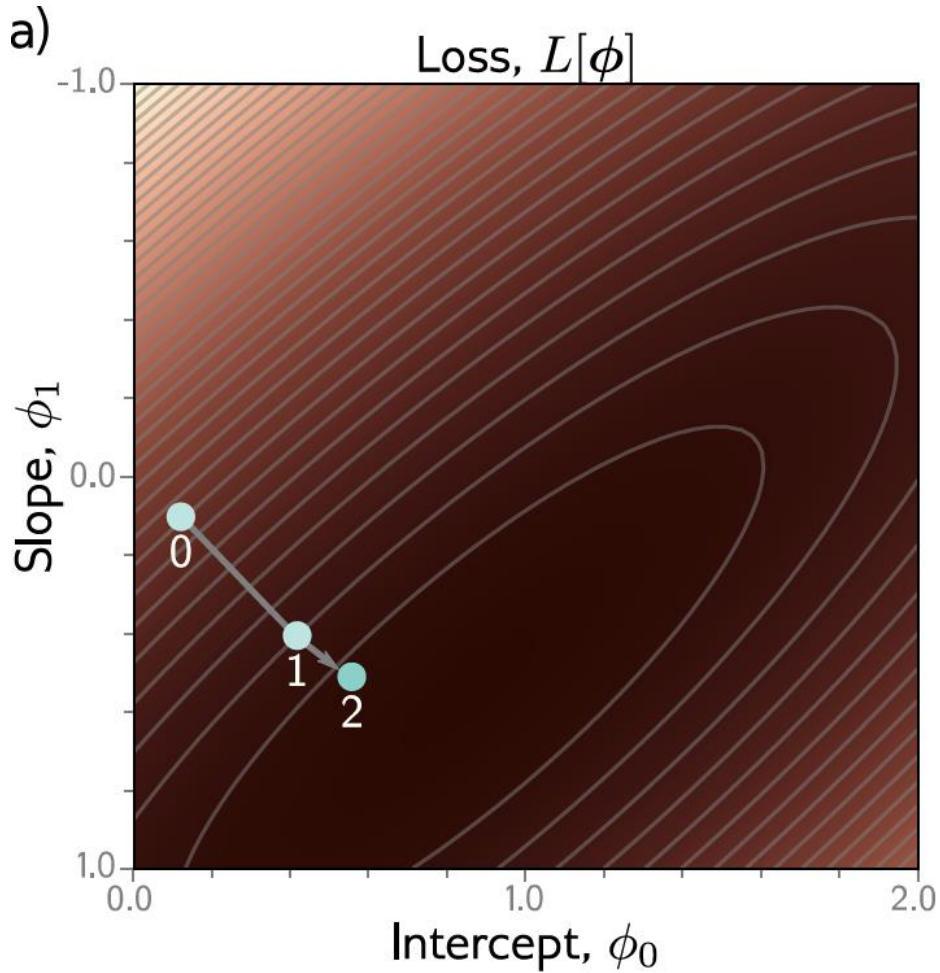
# Example: 1D Linear regression loss function



Loss function:

$$L[\phi] = \sum_{i=1}^{I} (\mathrm{f}[x_i, \phi] - y_i)^2$$

$$= \sum_{i=1}^{I} (\phi_0 + \phi_1 x_i - y_i)^2$$

"Least squares loss function"

# Example: 1D Linear regression loss function

# Example: 1D Linear regression loss function

# Example: 1D Linear regression loss function

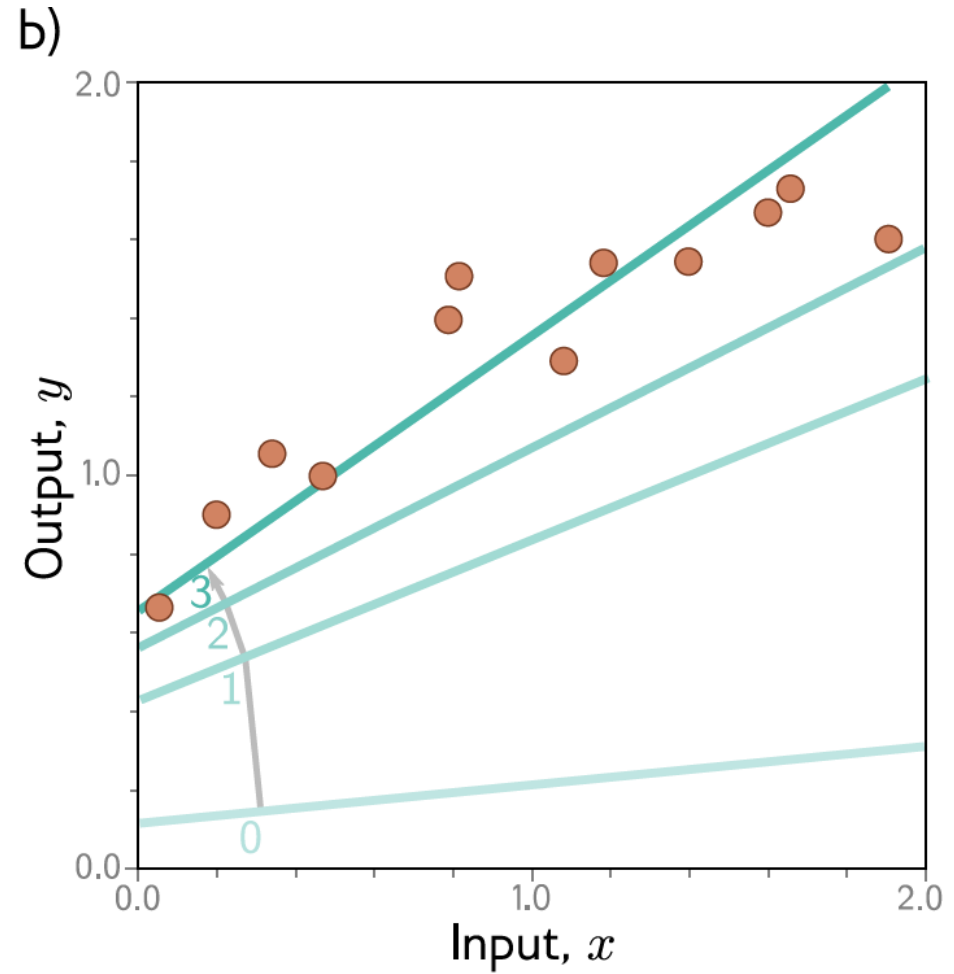# Example: 1D Linear regression loss function
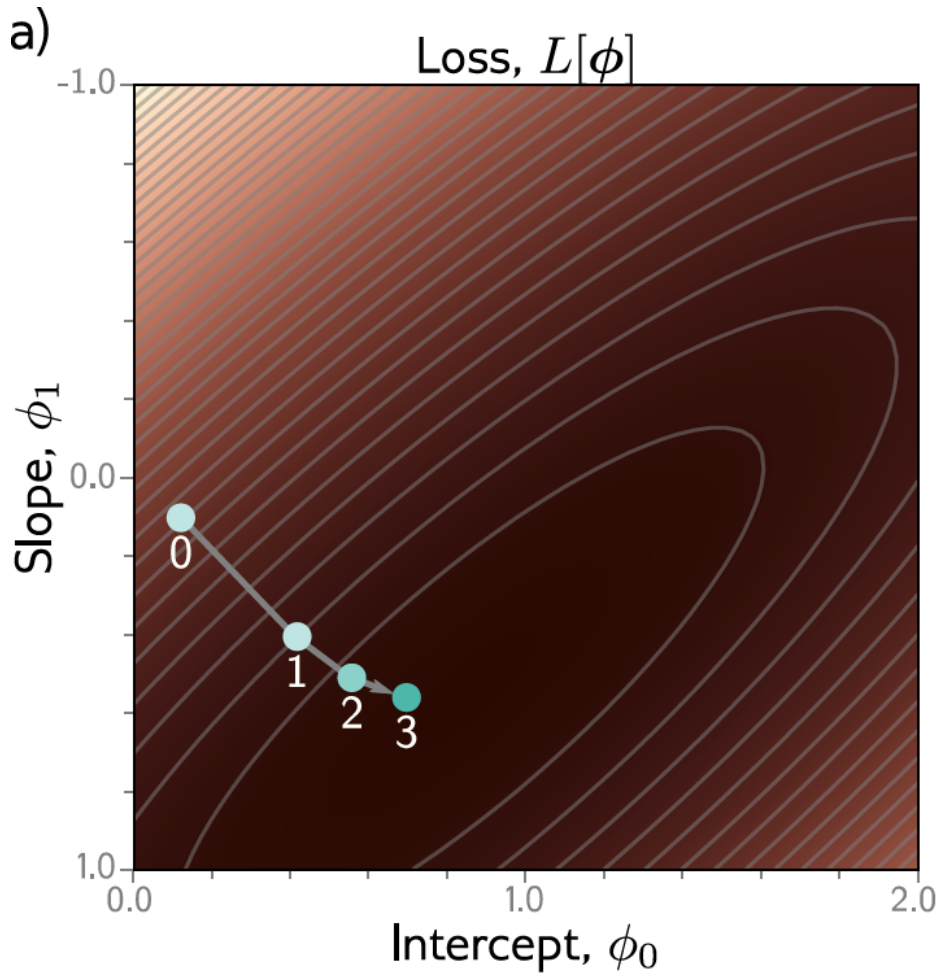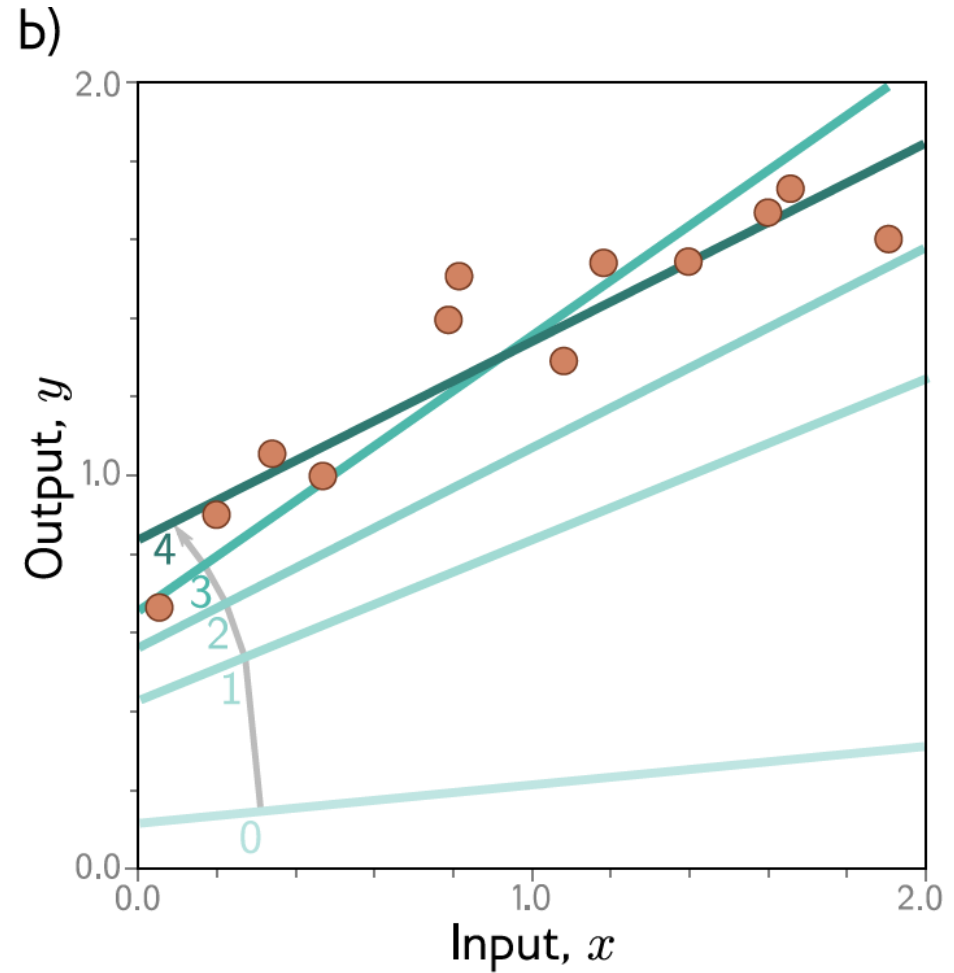


a)

b)

# Example: 1D Linear regression training
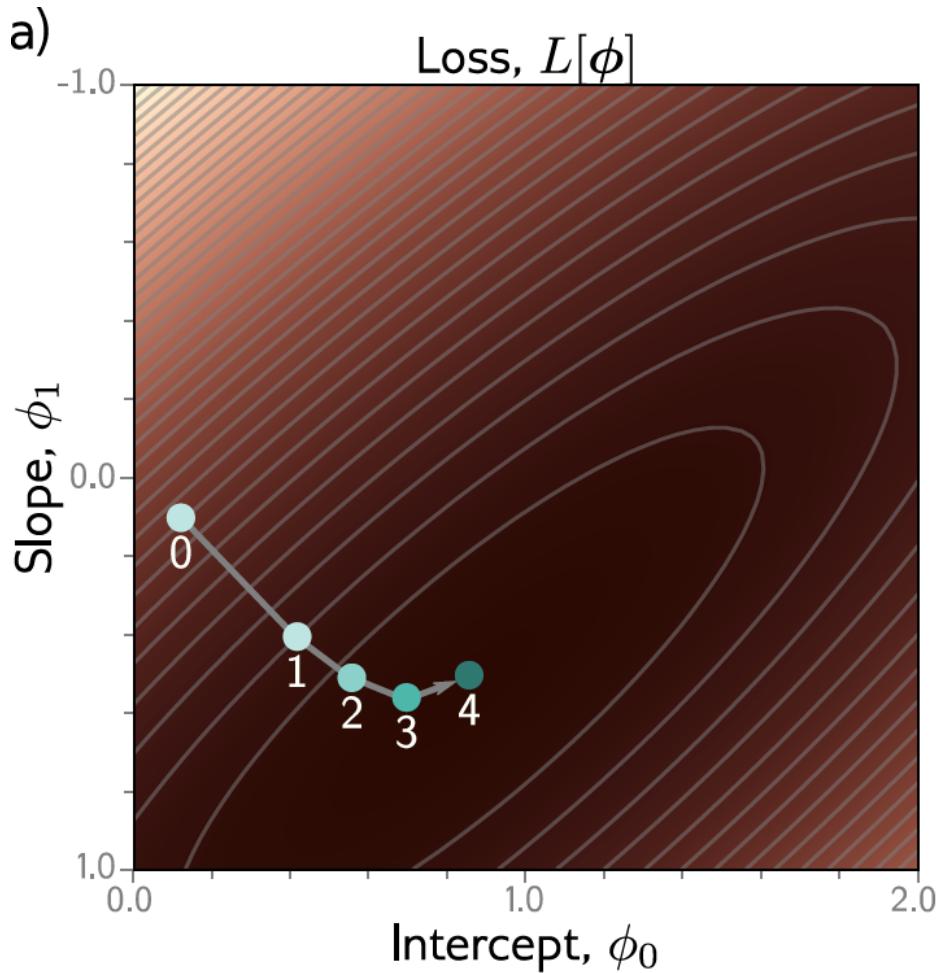
# Example: 1D Linear regression training

# Example: 1D Linear regression training

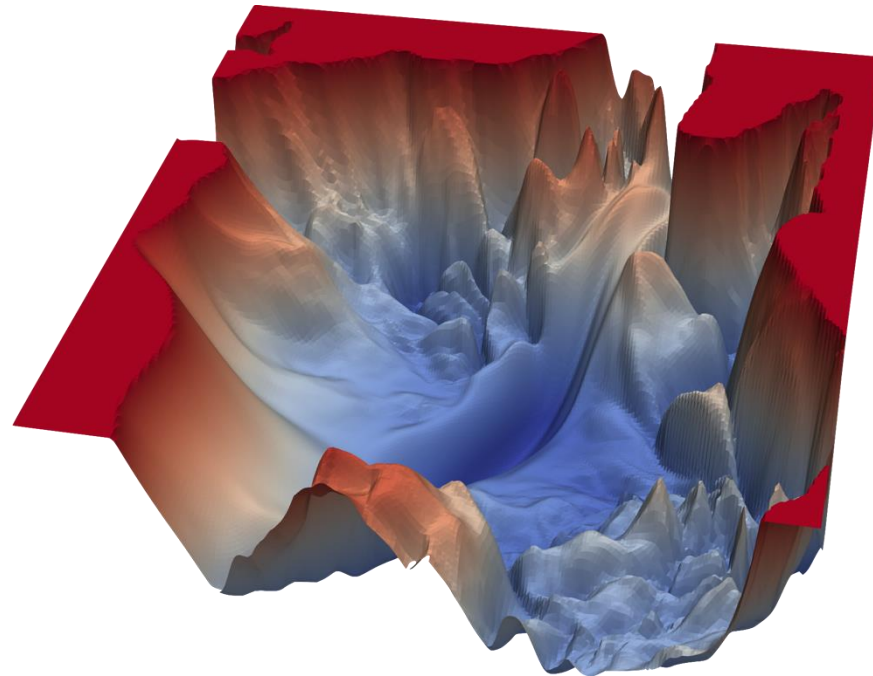# Example: 1D Linear regression training

# Example: 1D Linear regression training



a) Loss, $L[\phi]$ — Slope, $\phi_1$ vs Intercept, $\phi_0$

b) Output, $y$ vs Input, $x$

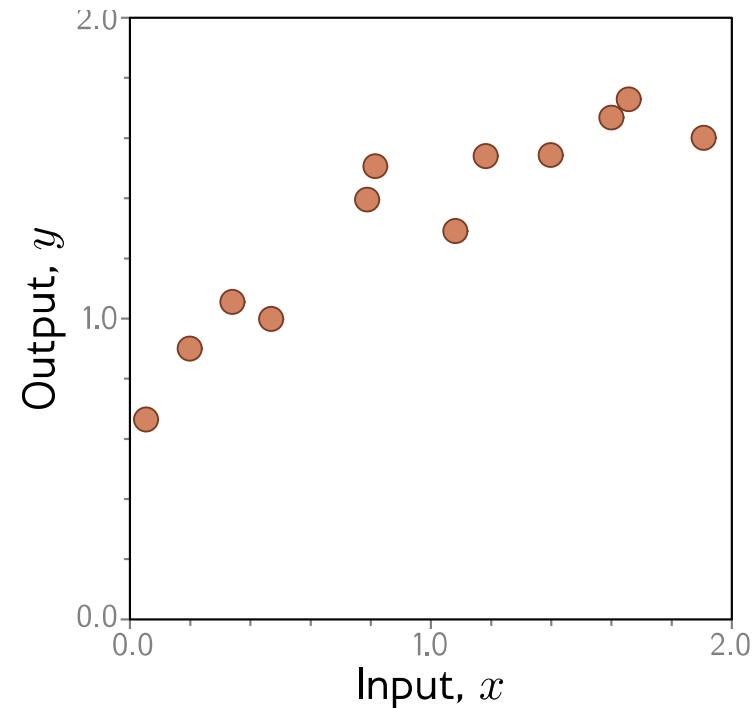This technique is known as gradient descent

# Possible objections

- But you can fit the line model in closed form!
  - Yes – but we won't be able to do this for more complex models

- But we could exhaustively try every slope and intercept combo!
  - Yes – but we won't be able to do this when there are a million parameters



Here's a visualization of the loss surface for the 56-layer neural network VGG-56(from
Visualizing the Loss Landscape of Neural Networks -- https://losslandscape.com/explorer

# Example: 1D Linear regression testing

- Test with different set of paired input/output data (Test Set)
    - Measure performance
    - Degree to which *Loss* is same as training = generalization

- Might not generalize well because
    - Model too simple: underfitting
    - Model too complex
        - fits to statistical peculiarities of data
        - this is known as overfitting

# Any Questions?

# Next Lecture

- How do we choose a loss function in a principled way?