

02-Linear Regression (Formulation)

Formulating a supervised learning task

- Variables related to the task can be separated into two parts
- x is the **input**, always known (during training and inference)
 - From now on, we assume x is a feature vector

$$x = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{pmatrix} \in \mathbb{R}^d, \quad d \text{ is the dimension}$$

- If a feature is real-valued, then it is just the value
- If a feature is categorical, then we may use **one-hot** representation. For example, if we would like to indicate if an animal is a cat, a dog, or an elephant, then we can use three binary indicators (although only one is one) to represent such information

whether a cat? $\rightarrow x_1$
a dog? $\rightarrow x_2$
an elephant? $\rightarrow x_3$

$$x_i \in \{0, 1\}$$

- Different features can be concatenated as a longer vector
- t is the **target/label/output**, which is of particular interest
Note: knowing x predicting $t \neq$ knowing t predicting x
 - In classification, t is a categorical variable
 - In regression, t is a real value, i.e.,

$$t \in \mathbb{R}$$

- For training, the machine learning model is given a set of data examples/samples in the form of (x, t) .

We name i th data sample $(x^{(i)}, t^{(i)})$.

Suppose we observe M samples in total for training, the **training set** is

$$\mathcal{D}_{\text{train}} = \{(x^{(i)}, t^{(i)})\}_{i=1}^M$$

- During inference, the model is asked to predict the label \hat{t}_* of an unseen input x_*
We may be given a set of unseen samples, and they are called a **test set**.

Formulating a supervised machine learning system

- Humans define a set of functions mapping from input to output, from which a machine learning system can pick up a particular function. This is known as a hypothesis class \mathcal{H} .

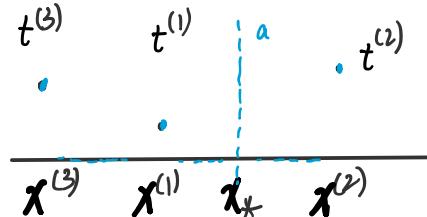
Usually, $\mathcal{H} \neq \{f: X \rightarrow Y\}$

Note: If we allow to learn from, or have no preference on, all functions, then learning is infeasible.

Suppose $x_* \notin \mathcal{D}_x$, consider a set of functions $\{f_a(x): a \in \mathbb{R}\}$, s.t.

$$f_a(x) = \begin{cases} y^{(i)}, & \text{if } x = x^{(i)} \\ a, & \text{if } x = x_* \\ 0, & \text{otherwise} \end{cases}$$

Since we only observe $\{(x^{(i)}, t^{(i)})\}_{i=1}^M$



and f_a performs equally well on the training set,

We have no preference on different values of a

The inference could be arbitrary.

Thus, we have to express our preference (hard and/or soft) on all possible functions, based on previous experience or simply for aesthetic purposes. Such preference comes with the machine learning system before training, and is called **inductive bias**.

It is also noted that the goal of machine learning systems is not to predict well on seen examples $x^{(1)}, x^{(2)}, \dots$, but rather on unseen examples x_*

- In the **training/learning phase**, the learning machine picks a particular function according to a certain criterion J , i.e., $h \in \mathcal{H}$

$$h = \arg \min_{h \in \mathcal{H}} J(h, \mathcal{D}_{\text{train}})$$

J is sometimes known as a **loss**, a **cost**, or a **training objective**.

"Statisticians seem to be pessimistic creatures who think in terms of losses. Decision theorists in economics and business talk instead in terms of gains (utility)."

--- James O. Berger (1985). *Statistical Decision Theory and Bayesian Analysis*.

- In the **inference/prediction phase**, the learning machine uses h to predict unseen samples \mathbf{x}_* , i.e.,

$$\hat{t}_* = h(\mathbf{x}_*)$$

In machine learning research, the model will be evaluated by a measure of success, which we call error

$$E(\underline{h(\mathbf{x}_*)}, t_*) \in \mathbb{R}$$

\hat{t}_*

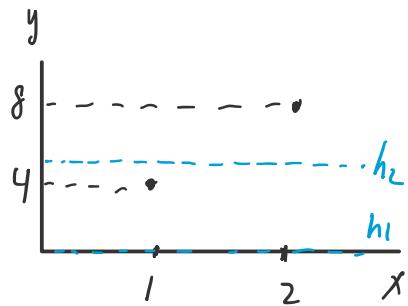
$$J(h) = \sum_{i=1}^n |t^{(i)} - h(x^{(i)})|^2$$

First toy example (finite hypothesis class)

$$x \in \mathbb{R}, t \in \mathbb{R} \quad \mathcal{D}_{\text{train}}: \begin{cases} (1, 4), (2, 8) \\ (x^{(1)}, t^{(1)}), (x^{(2)}, t^{(2)}) \end{cases}$$

Machine Learning System 1:

- $\mathcal{H} = \{h_1, h_2\}$
where $h_1(x) = 0, h_2(x) = 5$
- $J(h) = \sum_{m=1}^2 |t^{(m)} - h(x^{(m)})|^2$
- Optimization algorithm: by enumeration



Training:

$$\begin{aligned} J(h_1) &= 4^2 + 8^2 = 80 \\ J(h_2) &= 1^2 + 3^2 = 10 \end{aligned} \quad \Rightarrow \text{pick } h = h_2$$

Inference: $\forall \mathbf{x}_*$, the predicted $\hat{t}_* = h(\mathbf{x}_*) = 5$

Machine Learning System 2

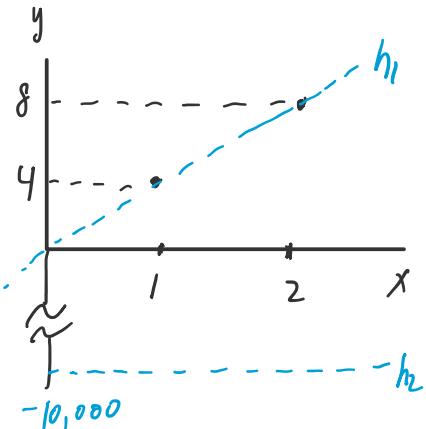
- $\mathcal{H} = \{h_1, h_2\}$

$$h_1(x) = 4x$$

$$h_2(x) = -10,000$$

- $J(h) = \sum_{n=1}^2 |h(x^{(n)}) + 10,000|^4$

- Optimization algorithm: by enumeration



Training: $J(h_1)$: very large } \Rightarrow pick $h = h_2$
 $J(h_2) = 0$

Inference:

$\forall x_*,$ the predicted label $\hat{t}_* = h_2(x) = -10,000$

Measure of Success:

Given test samples $x_*^{(1)}, t_*^{(1)} = (1.2, 4.5)$
 $x_*^{(2)}, t_*^{(2)} = (1.7, 7)$
 $x_*^{(3)}, t_*^{(3)} = (2.2, 8.5)$

The measure of success (error) is defined to be

$$E = \sum_{i=1}^3 |h(x_*^{(i)}) - t_*^{(i)}|$$

For Machine learning system 1: $E = 5.5$

For Machine learning system 2: $E = \dots$ (very large)

Conclusion: Machine learning system 1 wins

Lessons learned:

- A machine learning system involves many aspects, including the hypothesis class, the training objective function, and the optimization algorithm
- A machine learning model has no control on the test data and the measure of success
 - In particular, the measure of success could be different from the training objective

The Linear Regression Model

Hypothesis class

Let input $\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{pmatrix} \in \mathbb{R}^d$, output $t \in \mathbb{R}$

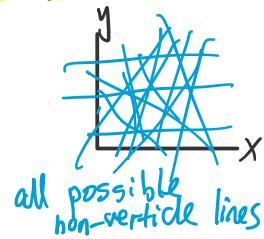
The hypothesis class is the affine functions on \mathbf{x} predicting t in the form of

$$h(\mathbf{x}) = w_1 x_1 + w_2 x_2 + \dots + w_d x_d + b = \sum_{i=1}^d w_i x_i + b$$

where $w = \begin{pmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{pmatrix} \in \mathbb{R}^d$ and $b \in \mathbb{R}$ are model parameters
called weights called bias

By linear algebra definitions, this can be denoted as

$$h(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$$



$$\text{The hypothesis class } \mathcal{H} = \{ h(\mathbf{x}) : h(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b, \mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R} \}$$

Note: linear vs affine transformation

Linear transformation takes the form of $h(\mathbf{x}) = \mathbf{w}^T \mathbf{x} = w_1 x_1 + \dots + w_d x_d$

Affine transformation allows the bias term

$$h(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$$

Intuitively, the bias term in linear regression learns an offset/baseline of the task.

Why "linear" regression?

- The simplest model
- We have to make assumptions (inductive bias)
- Non-linear problems can be transformed into a linear problem (usually in a higher dimensional space)
 - Human feature engineering
 - Kernel methods (implicit high dimensional space defined by inner-product)
 - Neural networks (composite non-linear functions)

Augmented feature

We define a trivial feature $x_0 = 1$ for all data points

We define the augmented feature representation

$$\tilde{\mathbf{x}} = \begin{bmatrix} x_0 \\ \mathbf{x} \end{bmatrix} = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{bmatrix} \in \mathbb{R}^{d+1}$$

Then, the linear regression model predicts

$$h(\tilde{\mathbf{x}}) = \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}$$

$$\text{where } \tilde{\mathbf{w}} := \begin{pmatrix} b \\ \mathbf{w} \end{pmatrix} = \begin{pmatrix} b \\ w_1 \\ \vdots \\ w_d \end{pmatrix} \in \mathbb{R}^{d+1}$$

We even drop " \sim " for simplicity.

Training objective of linear regression

- Suppose the parameters of a linear regression model is \mathbf{w}
[feature augmentation is implicit]
- We often define the loss as the square error
For a training sample $(\mathbf{x}^{(m)}, t^{(m)})$

$$\begin{aligned} \text{The per-sample loss is } J^{(m)} &= \frac{1}{2} (h(\mathbf{x}^{(m)}) - t^{(m)})^2 \\ &= \frac{1}{2} (\mathbf{w}^T \mathbf{x}^{(m)} - t^{(m)})^2 \end{aligned}$$

Then, the overall loss could be the average of per-sample loss

$$J = \frac{1}{M} \sum_{m=1}^M J^{(m)} = \frac{1}{2M} \sum_{m=1}^M (\mathbf{w}^T \mathbf{x}^{(m)} - t^{(m)})^2$$

Remarks:

- The square loss has profound theoretical justifications and we will explore some in future lectures. At this moment, we can intuitively think of it such that square error penalizes more on bad predictions
- $1/2$ is a constant having no effect in the optimization objective. It's included here mainly for aesthetic purposes (simplifying the derivative)
- Likewise, $1/M$ has not effect either, but in numerical optimization, including $1/M$ makes the model more stable with different sizes of data
- The model is also called ordinary least square (OLS).

Matrix representations

- The input can be represented as a matrix

$$X = \begin{bmatrix} \cdots & (\mathbf{x}^{(1)})^T \\ \cdots & (\mathbf{x}^{(2)})^T \\ \vdots & \\ \cdots & (\mathbf{x}^{(M)})^T \end{bmatrix}_{M \times (d+1)}$$

- The predicted targets are $\mathbf{y} = X \mathbf{w}$
 $M \times 1 \quad M \times (d+1) \quad (d+1) \times 1$
 (also denoted by $\hat{\mathbf{t}}$)

- The loss is

$$J = \frac{1}{2M} \|\mathbf{y} - \mathbf{t}\|_2^2$$

where $\|\cdot\|_2$ is the l_2 -norm of a vector

for $\mathbf{v} \in \mathbb{R}^d$

$$\|\mathbf{v}\|_2 = \sqrt{v_1^2 + v_2^2 + \cdots + v_d^2}$$

Or in general

$$\|\mathbf{v}\|_p = \left(\sum_{i=1}^d |v_i|^p \right)^{1/p}$$