



Institute for Advanced Studies
in Basic Sciences
Gava Zang, Zanjan, Iran

Regression with Linear Models

By Dr. Parvin Razzaghi (p.razzaghi@iasbs.ac.ir)

Winter 2024





Slides credit

- Professor:
 - Dr. Parvin Razzaghi
- University:
 - University of Toronto and Vector Institute
- Original Slides:
 - Amir-massoud Farahman
- Acknowledgment:
 - Credit for slides goes to many members of the ML Group at the U of T, and beyond, including (recent past): Roger Grosse, Murat Erdogdu, Richard Zemel, Juan Felipe Carrasquilla, Emad Andrews. For any comments or suggestions, please contact: p.razzaghi@iasbs.ac.ir

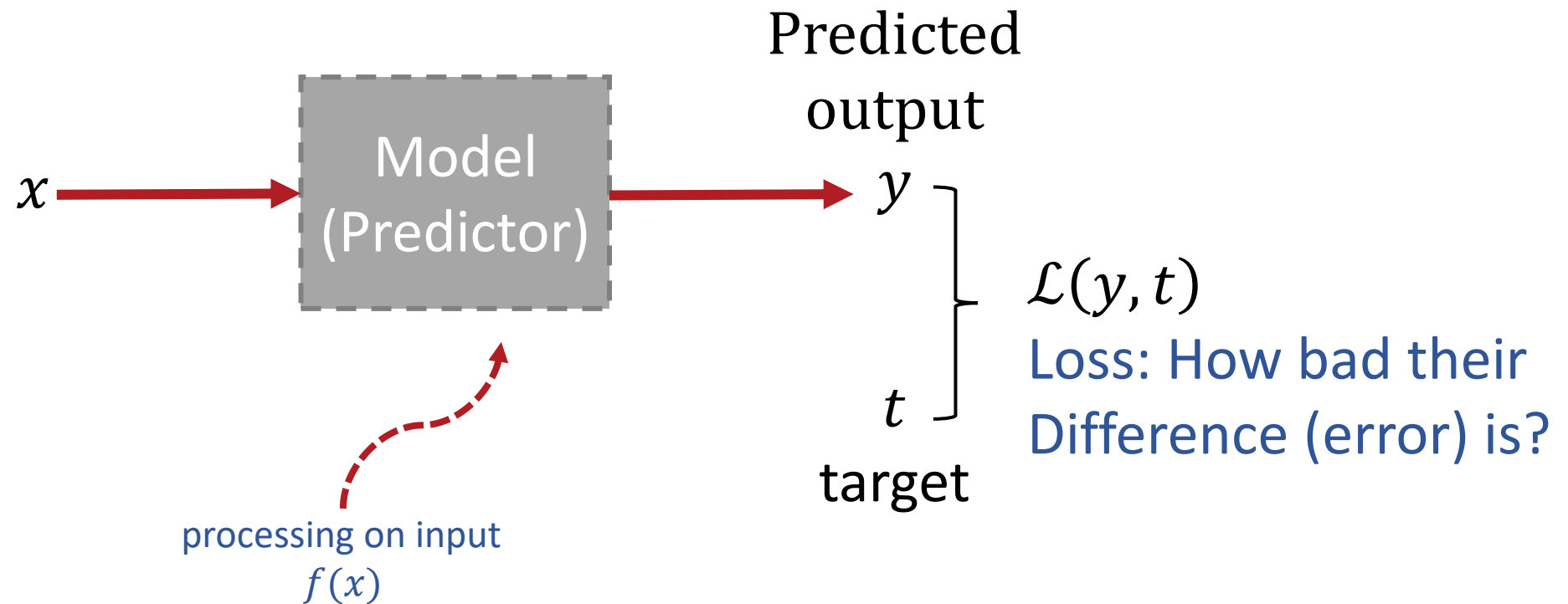


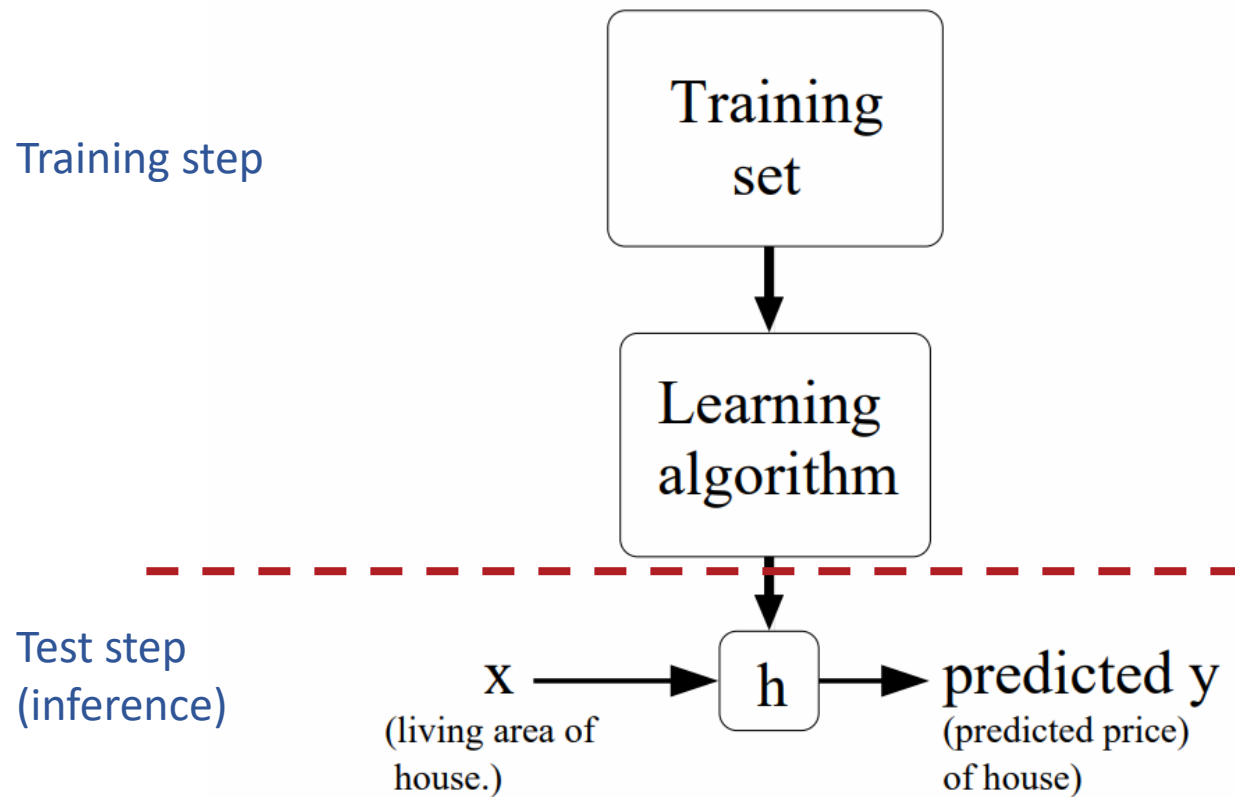
Table of Contents

- Modular Approach to ML
- Regression
 - Linear Regression
 - Basic Expansion
 - Regularization
 - Probabilistic Interpretation of the Squared Error
 - Bias-variance tradeoff
 - Bayesian linear regression



Modular Approach ML Algorithm Design







Modular Approach ML Algorithm Design

- we will take a more **modular** approach:
 - define the problem
 - provide a training dataset
 - choose a **model** describing the relationships between variables of interest
 - define a **loss function** quantifying how bad the fit to the data is
 - (possibly) choose a **regularizer** saying how much we prefer different candidate models (or explanations of data), **before (prior to)** seeing the data
 - fit the model that minimizes the loss function and satisfy the constraint/penalty imposed by the regularizer, possibly using an **optimization algorithm**
- Mixing and matching these modular components gives us a lot of new ML methods.



Define the problem & provide a dataset



The Supervised Learning Setup

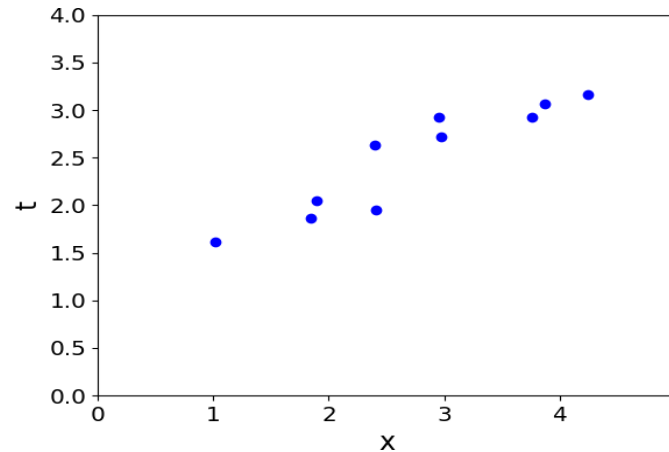
- Features
 - Living area
 - Number of bedrooms
 - ...
- Target variable
 - Price



Living area (feet ²)	#bedrooms	Price (1000\$s)
2104	3	400
1600	3	330
2400	3	369
1416	2	232
3000	4	540
⋮	⋮	⋮



The Supervised Learning Setup



Recall that in supervised learning:

- There is a target $t \in T$ (also called response, outcome, output, class)
- There are features $\mathbf{x} \in X$ (also called inputs or covariates)

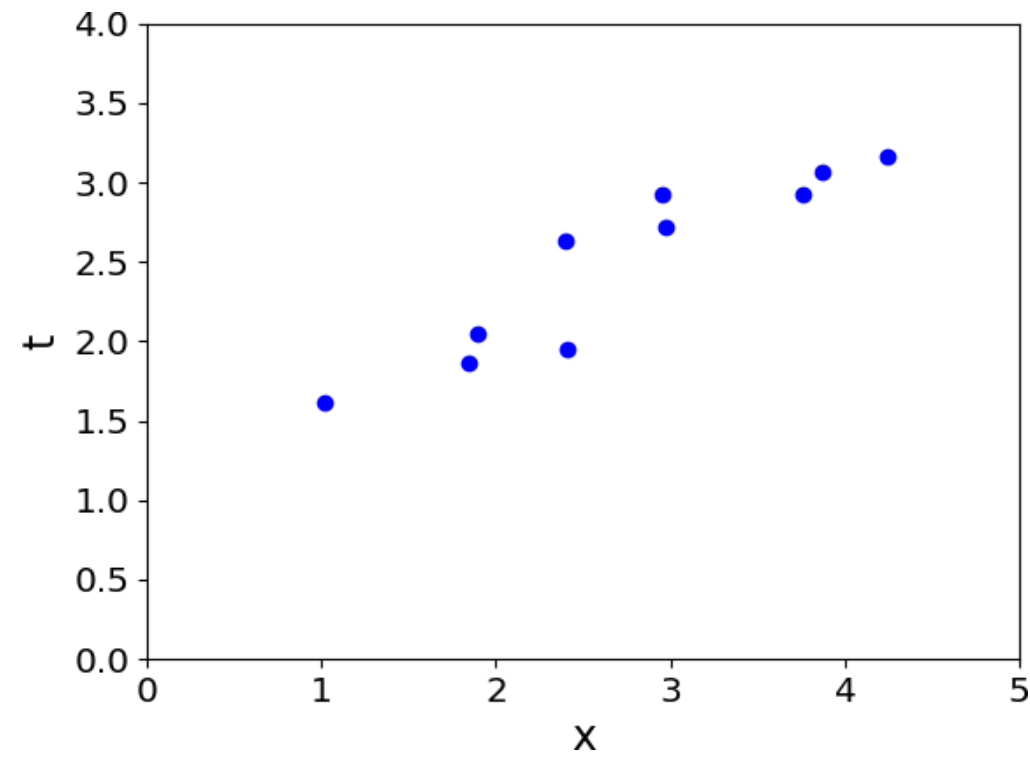
The goal is to learn a function $f : X \rightarrow T$ such that

$$t \approx y = f(\mathbf{x}),$$

based on given data $D = \{(\mathbf{x}^{(i)}, t^{(i)}) \text{ for } i = 1, 2, \dots, N\}$.

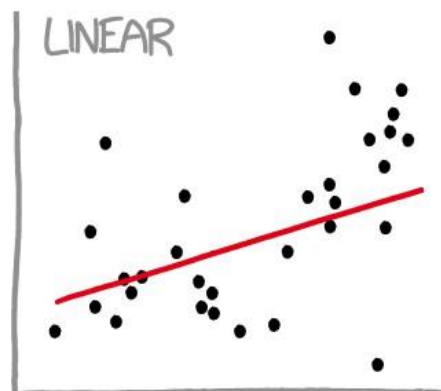


choose a **model** describing the relationships between variables of interest?





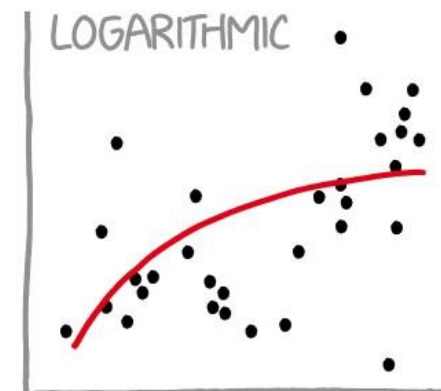
Regression with Linear Models



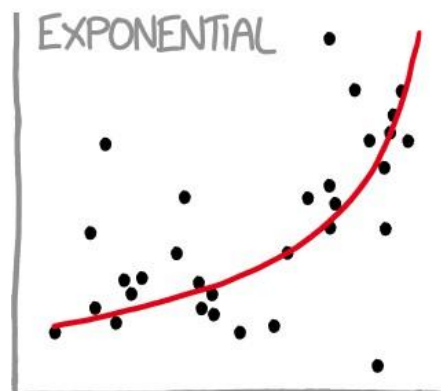
"HEY, I DID A
REGRESSION."



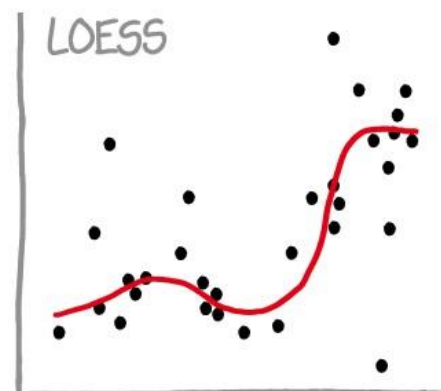
"I WANTED A CURVED
LINE, SO I MADE ONE
WITH MATH."



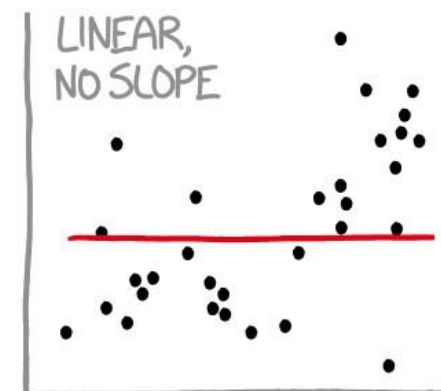
"LOOK, IT'S
TAPERING OFF!"



"LOOK, IT'S GROWING
UNCONTROLLABLY!"



"I'M SOPHISTICATED, NOT
LIKE THOSE BUMBLING
POLYNOMIAL PEOPLE."



"I'M MAKING A
SCATTER PLOT BUT
I DON'T WANT TO."



Linear Regression – Model

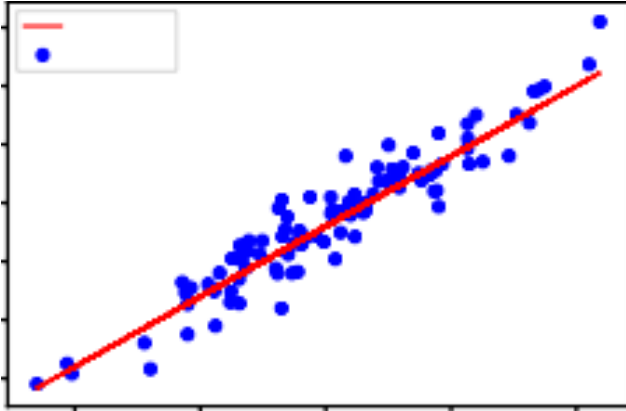
- **Model:** In linear regression, we use linear functions of the inputs $\mathbf{x} = (x_1, \dots, x_D)$ to make predictions y of the target value t :

$$y = f(\mathbf{x}) = \sum_j w_j x_j + b$$

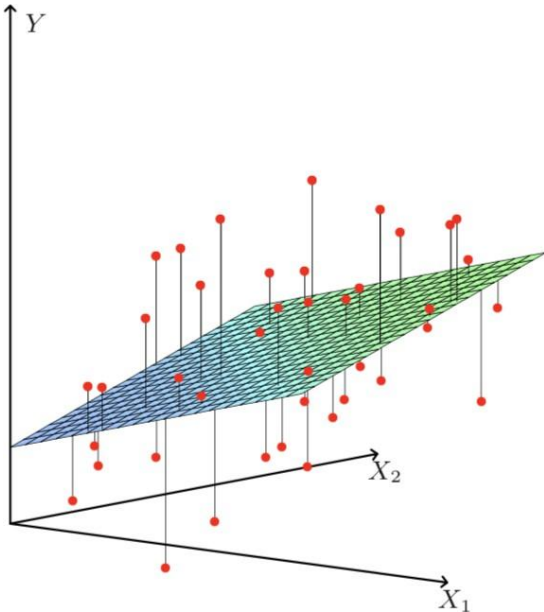
- y is the **prediction**
- \mathbf{w} is the **weights**
- b is the **bias** (or **intercept**) (do not confuse with the bias-variance tradeoff in the next lecture)
- \mathbf{w} and b together are the **parameters**
- We hope that our prediction is close to the target: $y \approx t$.



What is Linear? 1 Feature vs. D Features



- If we have only 1 feature:
 $y = wx + b$ where $w, x, b \in \mathbb{R}$
- y is linear in x

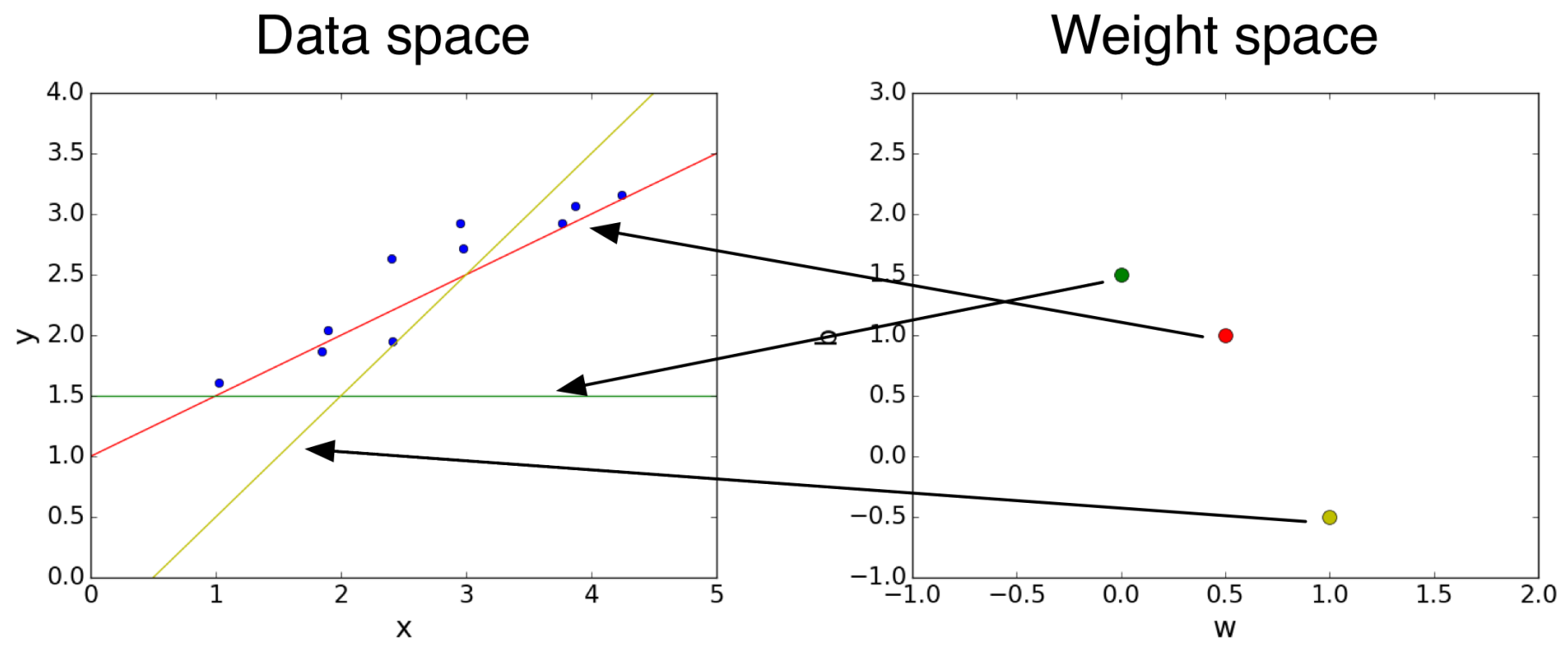


- If we have D features:
 $y = \mathbf{w}^T \mathbf{x} + b$ where $\mathbf{w}, \mathbf{x} \in \mathbb{R}^D, b \in \mathbb{R}$
- y is linear in \mathbf{x}

Relation between the prediction y and inputs \mathbf{x} is linear in both cases.



Weight Space vs. Data Space



Recall that:

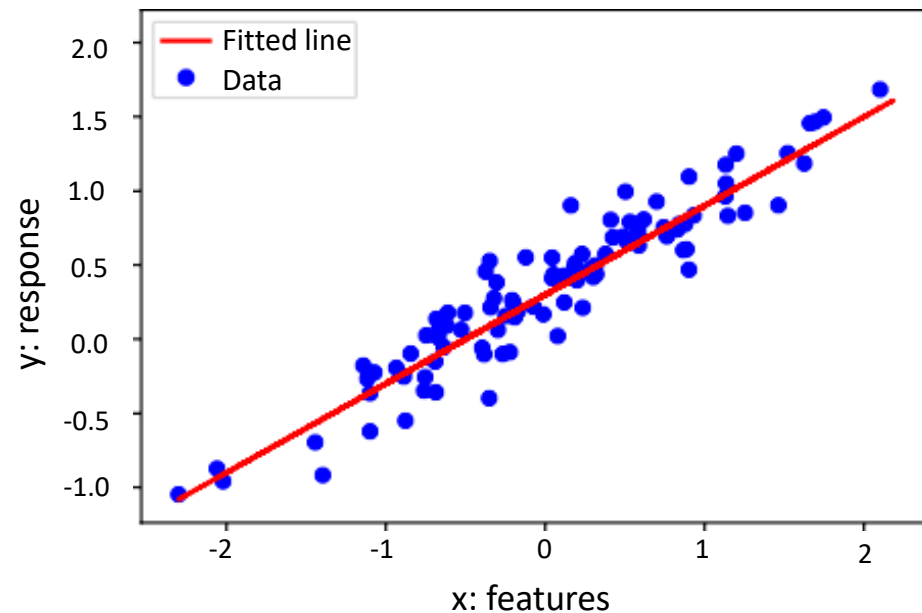
$$y = f(x) = \sum_j w_j x_j + b$$



Linear Regression

We have a dataset $D = \{(x^{(i)}, t^{(i)})\}_{i=1}^N$ where:

- $x^{(i)} = (x_1^{(i)}, x_2^{(i)}, \dots, x_D^{(i)})^T \in \mathbb{R}^D$ are the inputs, e.g., age, height
- $t^{(i)} \in \mathbb{R}$ is the target or response, e.g., income
- predict $t^{(i)}$ with a linear function of $x^{(i)}$:



- $t^{(i)} \approx y^{(i)} = \mathbf{w}^T \mathbf{x}^{(i)} + b$
- Find the “best” line (w, b) .
- **Q:** How should we find the **best** line?

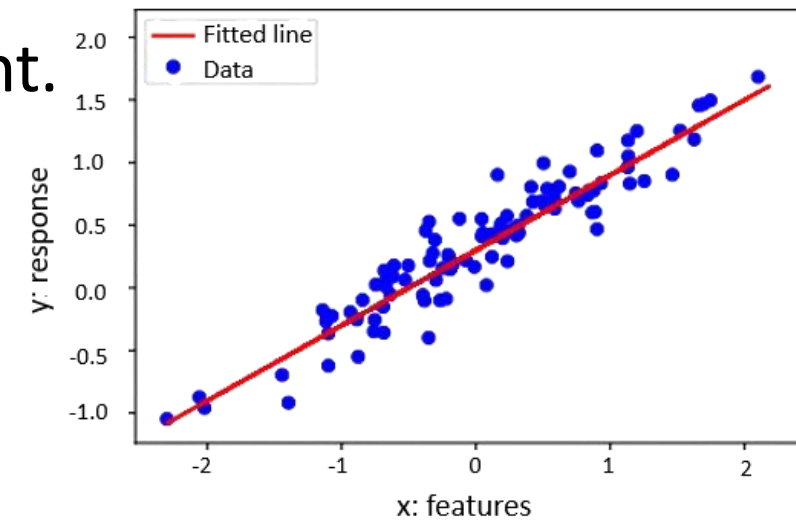


Linear Regression – Loss Function

- How to quantify the quality of the fit to data?
- A **loss function** $\mathcal{L}(y, t)$ defines how bad it is if, for some input x , the algorithm predicts y , but the target is actually t .
- **Squared error loss function:**

$$\mathcal{L}(y, t) = \frac{1}{2} (y - t)^2$$

- $y - t$ is the **residual**, and we want to make its magnitude small
- The $\frac{1}{2}$ factor is just to make the calculations convenient.





Linear Regression – Loss Function

- **Cost function:** loss function averaged over all training examples

$$\mathcal{J}(\mathbf{w}, b) = \frac{1}{N} \sum_{i=1}^N \mathcal{L}(y^{(i)}, t^{(i)})$$

$$= \frac{1}{2N} \sum_{i=1}^N (y^{(i)} - t^{(i)})^2$$

$$= \frac{1}{2N} \sum_{i=1}^N (\mathbf{w}^T x^{(i)} + b - t^{(i)})^2$$



Linear Regression – Loss Function

- To find the best fit, we find a model (parameterized by its weights \mathbf{w} and b) that minimizes the cost:

$$\underset{(\mathbf{w}, b)}{\text{minimize}} \mathcal{J}(\mathbf{w}, b) = \frac{1}{N} \sum_{i=1}^N \mathcal{L}(y^{(i)}, t^{(i)})$$

- The terminology is not universal. Some might call “loss” **pointwise loss** and the “cost function” the **empirical loss** or **average loss**.



Vector Notation

- We can organize all the training examples into a **design matrix \mathbf{X}** with one row per training example, and all the targets into the **target vector \mathbf{t}** .

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}^{(1)T} \\ \mathbf{x}^{(2)T} \\ \mathbf{x}^{(3)T} \end{bmatrix} = \begin{bmatrix} 8 & 0 & 3 & 0 \\ 6 & -1 & 5 & 3 \\ 2 & 5 & -2 & 8 \end{bmatrix}$$

One feature across
all training examples

One training
example (vector)

- Computing the predictions for the whole dataset:

$$\mathbf{X}\mathbf{w} + b\mathbf{1} = \begin{bmatrix} \mathbf{w}^T \mathbf{x}^{(1)} + b \\ \vdots \\ \mathbf{w}^T \mathbf{x}^{(N)} + b \end{bmatrix} = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(N)} \end{bmatrix} = \mathbf{y}$$



Vectorization

- Computing the squared error cost across the whole dataset:

$$\mathbf{y} = \mathbf{X}\mathbf{w} + b\mathbf{1}$$

$$\mathcal{J} = \frac{1}{2N} \|\mathbf{y} - \mathbf{t}\|^2$$

- Note that sometimes we may use $\mathcal{J} = \frac{1}{2N} \|\mathbf{y} - \mathbf{t}\|^2$, without $\frac{1}{N}$ normalizer. That would correspond to the sum of losses, and not the average loss. That does not matter as the minimizer does not depend on N .



Vectorization

- We can also add a column of 1s to the design matrix, combine the bias and the weights, and conveniently write

$$X = \begin{bmatrix} 1 & [x^{(1)}]^T \\ 1 & [x^{(2)}]^T \\ 1 & \vdots \end{bmatrix} \in \mathbb{R}^{N \times D+1} \text{ and } w = \begin{bmatrix} b \\ w_1 \\ w_2 \\ \vdots \end{bmatrix} \in \mathbb{R}^{D+1}$$

Then, our predictions reduce to $\mathbf{y} = \mathbf{X}\mathbf{w}$.

Any questions?

