

Lecture 01:

“Supervised Learning”

Supervised Learning

Training Input Training Output

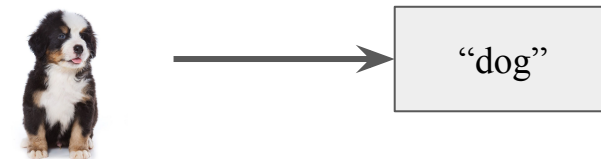
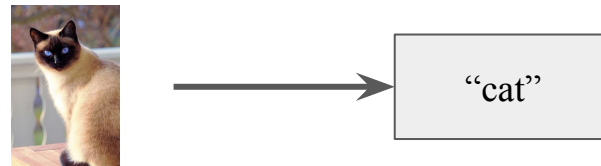


...

Test Input Test Output

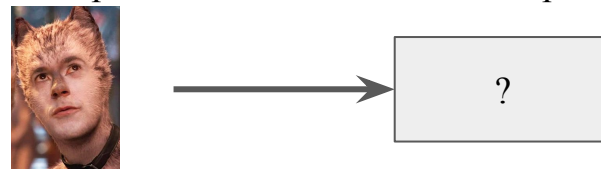


Training Input Training Output



...

Test Input Test Output



3 steps:

1

Choose a form for the function which relates inputs (x) to output (y)

2

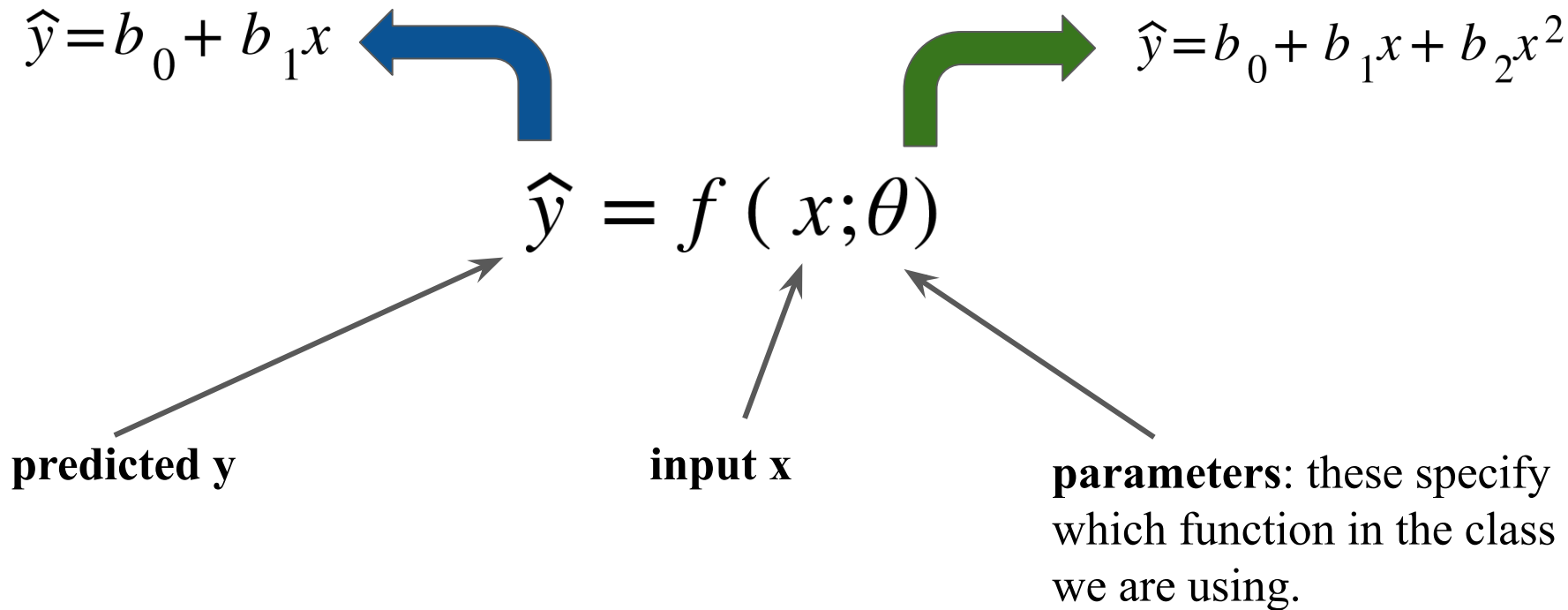
Define training loss

3

Find function in a form which gives the smallest training loss

1

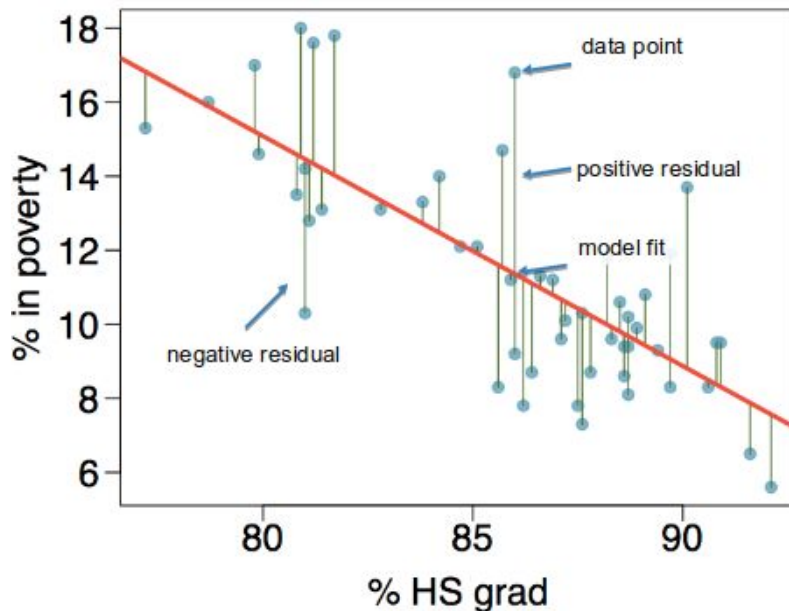
Choose a form for the function which relates inputs (x) to output (y)



2

Define training loss

- A **loss function** measures the deviation of a model's fit from observed data
- We minimize our loss function by tweaking the parameters, given training data



- **Residuals** are the errors from the model fit
- $\text{Data} = \text{Fit} + \text{Residual}$

We want a line with small residuals...

Option 1: Minimize the sum of magnitudes (absolute values) of residuals: **The L_1 -norm** (also called LAD or Least Absolute Deviation)

$$L(\theta) = \sum_{i=1}^n |y_i - \hat{y}_i| = \sum_{i=1}^n |r_i| = \|\mathbf{r}\|_1$$

Option 2: Minimize the sum of squared residuals: **The squared L_2 -norm** (also called OLS or Ordinary Least Squares)

$$L(\theta) = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n r_i^2 = \|\mathbf{r}\|_2^2$$

The most commonly used is **least squares**

- Solutions can be easily computed
- Big errors count relatively more than small errors


3

Find function in a form which gives the smallest training loss

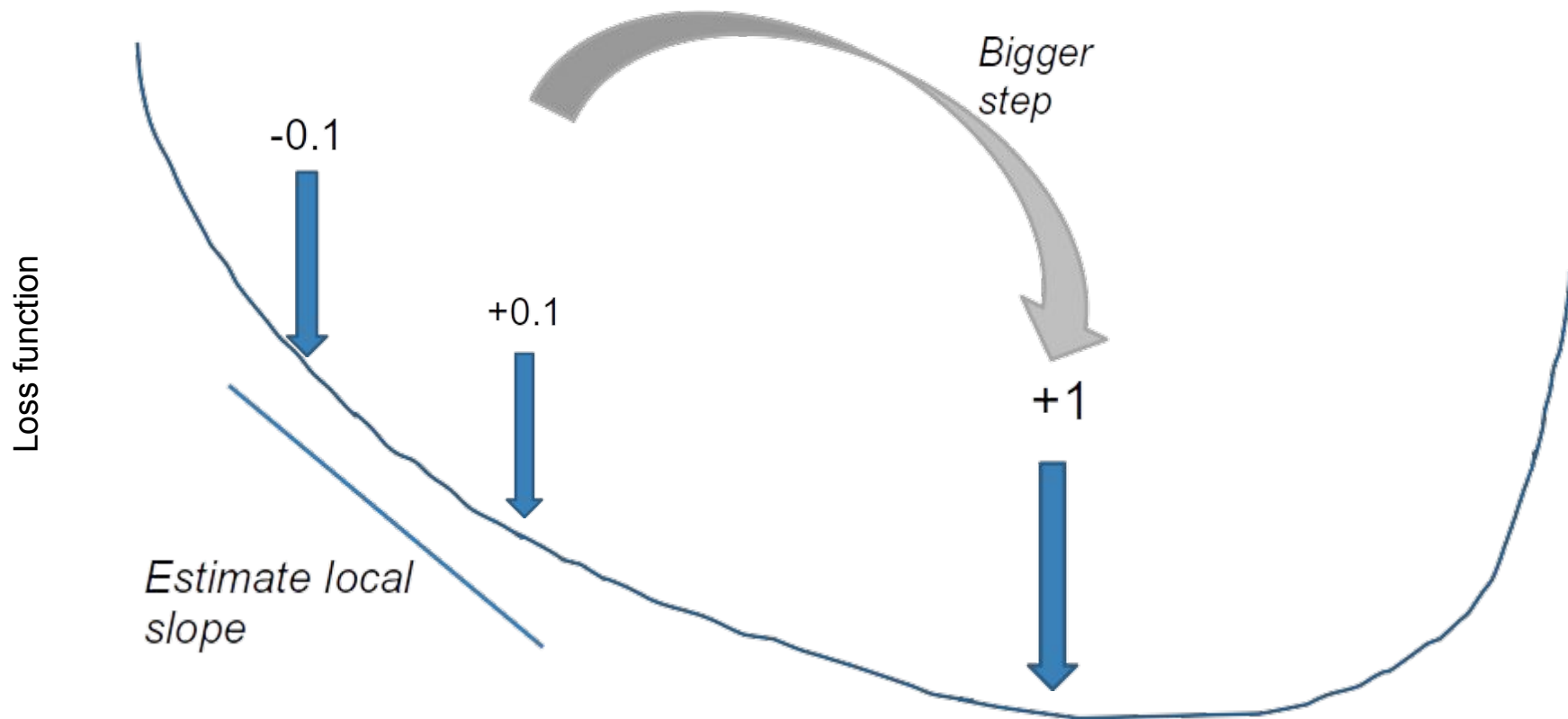
- This is equivalent to finding the **best parameters**.
- We want to minimize the training loss by trying different parameter values.

$$y = f(x; \theta) \xrightarrow{\text{red arrow}} \hat{y} = b_0 + b_1 x$$

$$b_0 = ?, b_1 = ? \xleftarrow{\text{blue arrow}} L(b_0, b_1) = ?$$



Minimizing the cost



Derivative of Loss

Why? To speed up optimization (instead of guess and iterate)

$$L = \sum_{i=1}^n \left(y_i - b_0 - b_1 x_i \right)^2$$

$$\begin{aligned} \frac{\partial L}{\partial b_0} &= -2 \sum_{i=1}^n \left(y_i - b_0 - b_1 x_i \right) \\ &= -2 \sum_{i=1}^n r_i \end{aligned}$$

$$\begin{aligned} \frac{\partial L}{\partial b_1} &= -2 \sum_{i=1}^n \left(y_i - b_0 - b_1 x_i \right) x_i \\ &= -2 \sum_{i=1}^n r_i x_i \end{aligned}$$

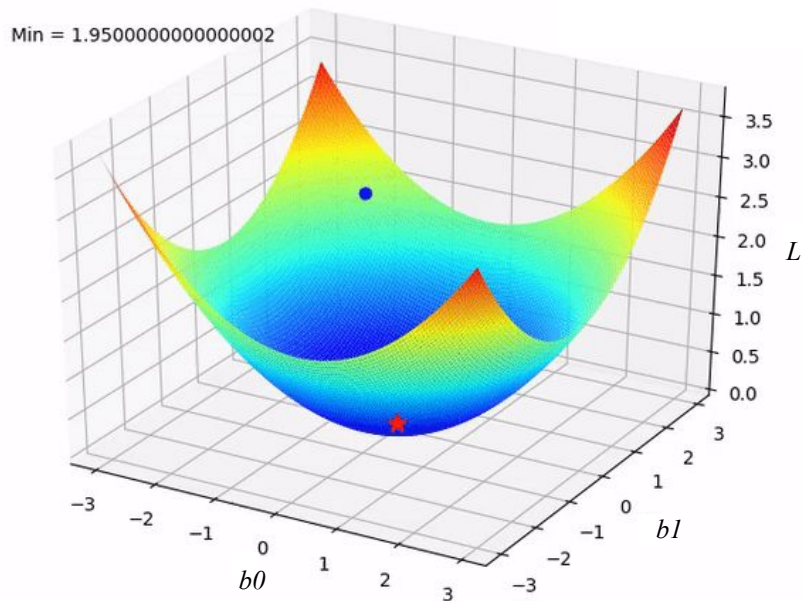
Derivative of Loss (L_1)

$$L = \sum_{i=1}^n |y_i - b_0 - b_1 x_i|$$

$$\begin{aligned} \frac{\partial L}{\partial b_0} &= - \sum_{i=1}^n \text{sgn}(y_i - b_0 - b_1 x_i) \\ &= - \sum_{i=1}^n \text{sgn}(r_i) \end{aligned}$$

$$\begin{aligned} \frac{\partial L}{\partial b_1} &= - \sum_{i=1}^n \text{sgn}(y_i - b_0 - b_1 x_i) \cdot x_i \\ &= - \sum_{i=1}^n \text{sgn}(r_i) \cdot x_i \end{aligned}$$

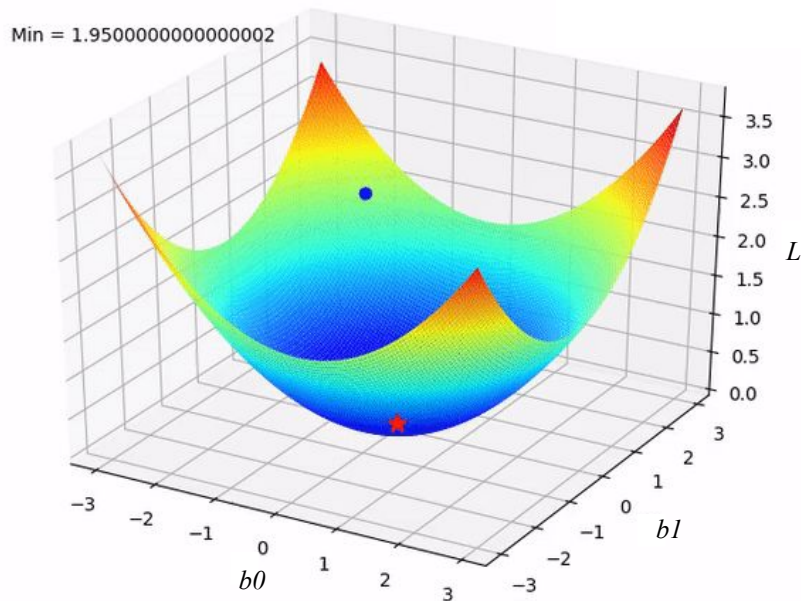
Derivative of loss (vectors)



$$\nabla_b L = \begin{bmatrix} \frac{\partial L}{\partial b_0} \\ \frac{\partial L}{\partial b_1} \end{bmatrix} = \begin{bmatrix} -2 \sum_{i=1}^n r_i \\ -2 \sum_{i=1}^n r_i x_i \end{bmatrix}$$

This is called the **Gradient** or **Jacobian**

Derivative of loss (vectors)



In vector notation...

$$\hat{\mathbf{y}} = \mathbf{X}\mathbf{b} \rightarrow \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \dots \\ \hat{y}_n \end{bmatrix} = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \dots & \dots \\ 1 & x_n \end{bmatrix} \begin{bmatrix} b_0 \\ b_1 \end{bmatrix}$$

(Residual) $\mathbf{r} = \mathbf{y} - \hat{\mathbf{y}}$

(Loss) $L = (\mathbf{y} - \mathbf{X}\mathbf{b})^T (\mathbf{y} - \mathbf{X}\mathbf{b})$

(Gradient) $\nabla_{\mathbf{b}} J = -2\mathbf{X}^T \mathbf{r}$

How good is the fit?

For linear regression we often use R^2 : the **coefficient of determination**

$$R^2 = 1 - \frac{RSS}{TSS}$$

Residual Sum of Squares

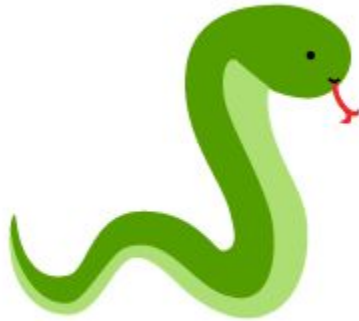
Total Sum of Squares

$$= 1 - \frac{\sum_{i=1}^n (y - \hat{y})^2}{\sum_{i=1}^n (y - \bar{y})^2}$$

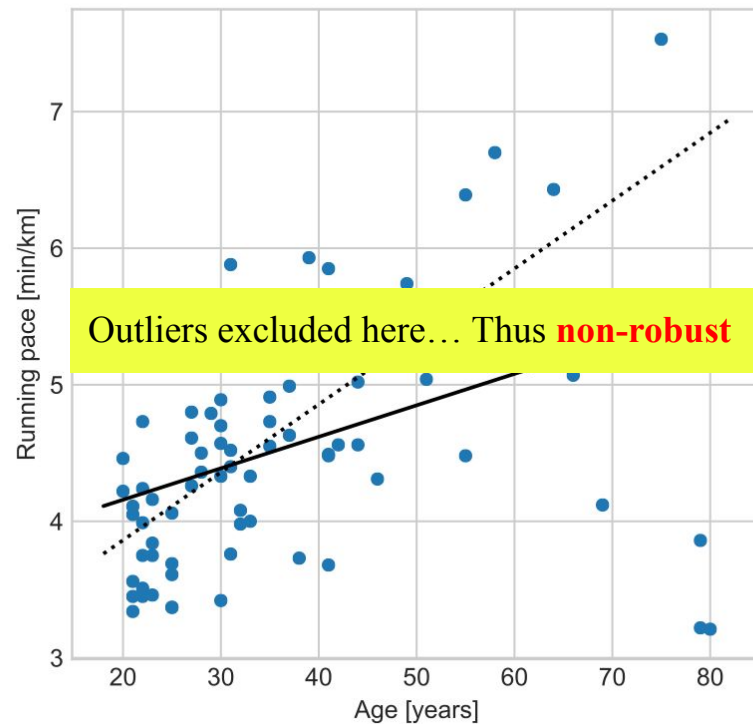
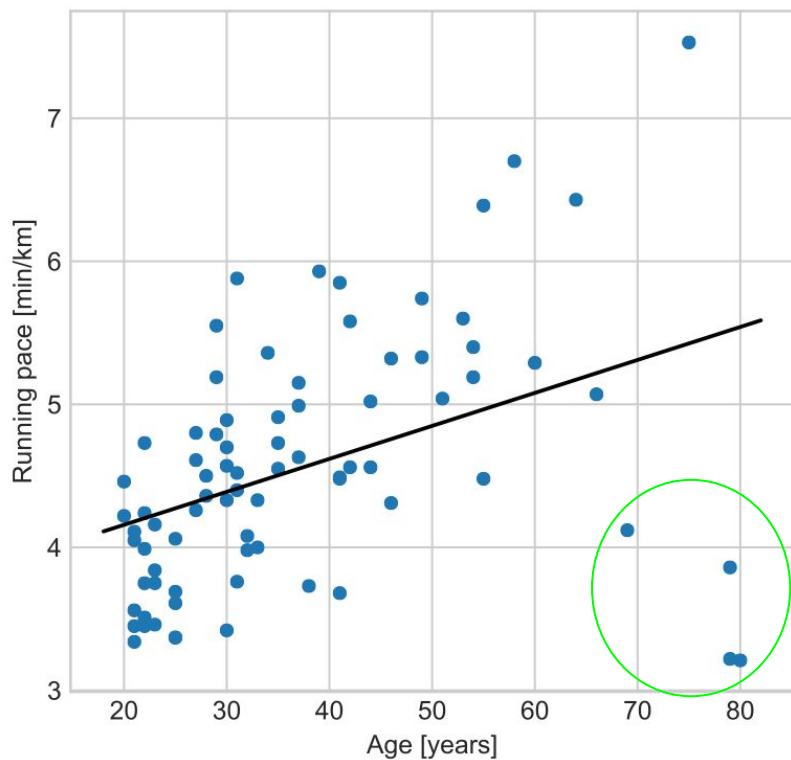
Mean value of the dataset

$R^2 = 0 \rightarrow$ no fit; $R^2 = 1 \rightarrow$ perfect fit. Note: OLS will always have the highest possible R^2 value.

Let's try it in Python...



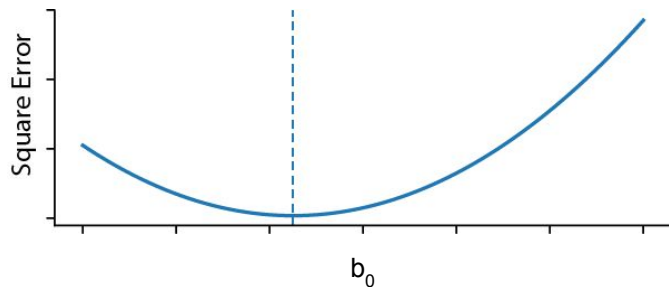
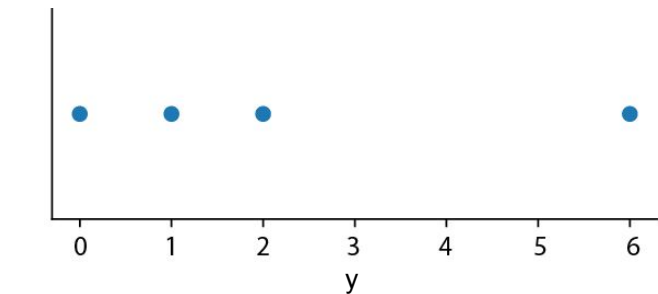
A Different Loss Function



Robust regression

Robust statistics: provides methods which are not affected by (less sensitive to) outliers

- **Mean** (a measure of central tendency) is **sensitive** to outliers, **median** is **robust** to outliers



Where is the sum of squared error minimal?

Recall:

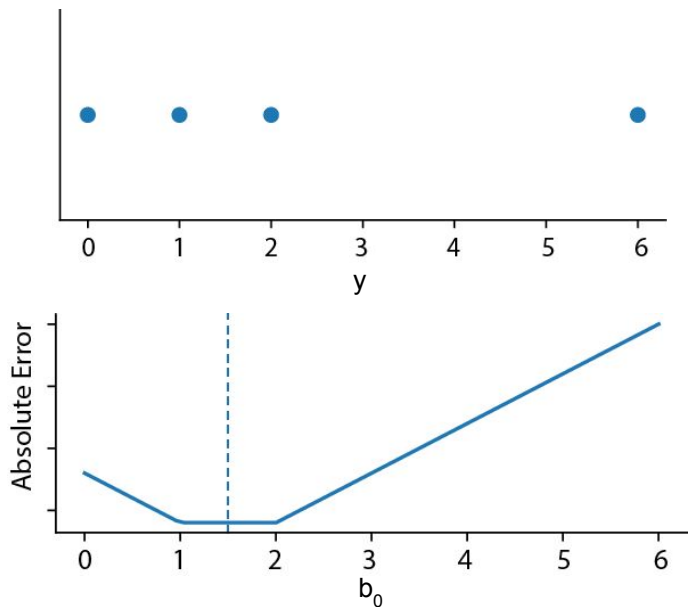
$$L = \sum_{i=1}^n (y_i - b_0)^2$$

$$\frac{\partial L}{\partial b_0} = -2 \sum_{i=1}^n (y_i - b_0)$$
$$\frac{\partial L}{\partial b_0} = 0 \rightarrow b_0 = \frac{\sum_{i=1}^n y_i}{n}$$

Answer: at the mean.

Robust regression

Question: Can we change the loss function to **improve its robustness to outliers**? Let's consider minimizing the sum of absolute errors.



Where is the sum of absolute errors minimal?

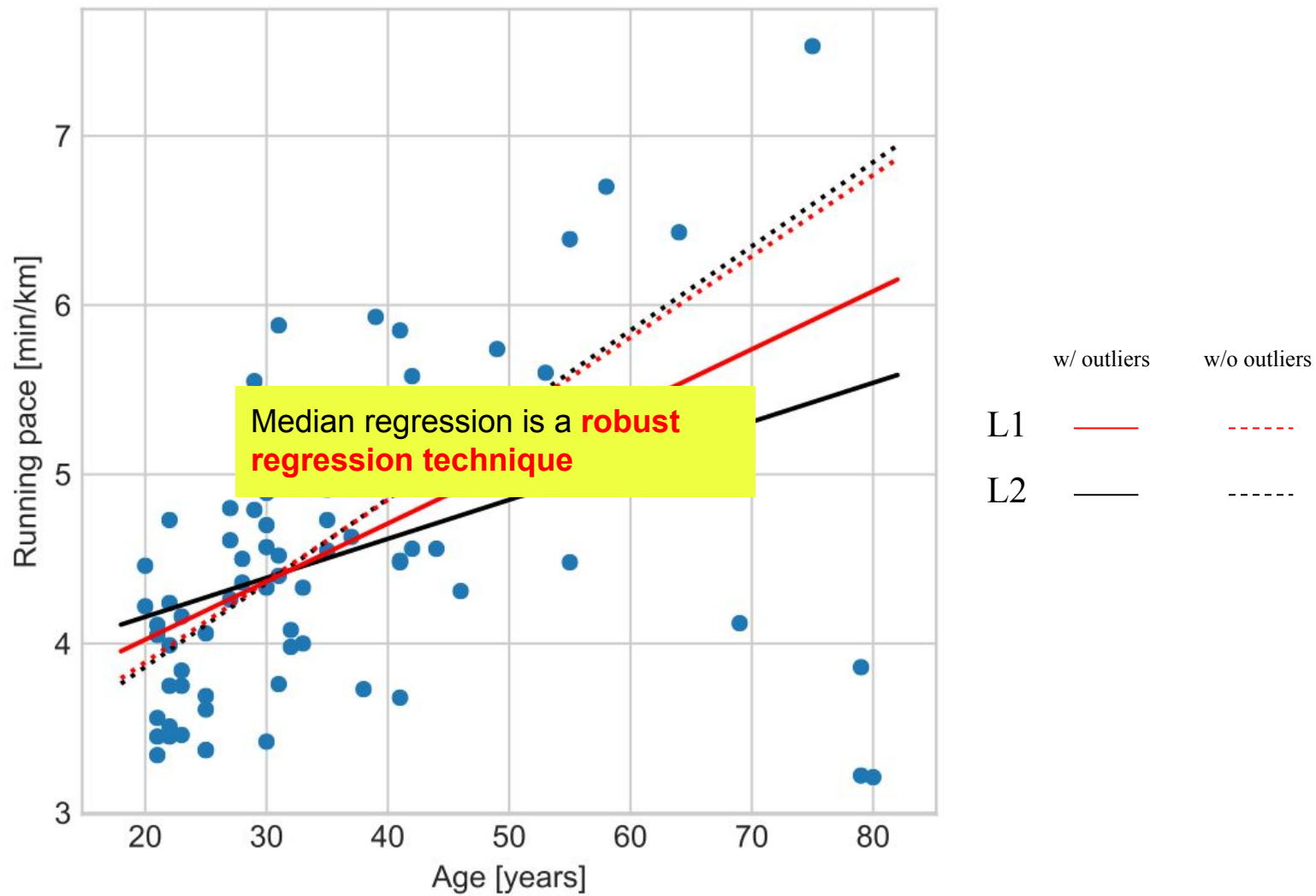
Recall:

$$L = \sum_{i=1}^n |y_i - b_0|$$

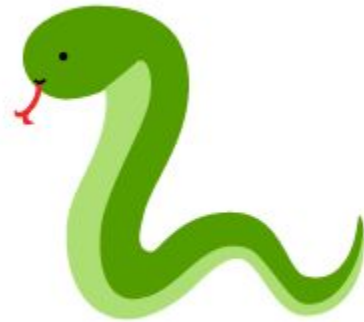
$$L = \sum_{i=1}^n \begin{cases} y_i - b_0 & \text{if } y_i > b_0 \\ - (y_i - b_0) & \text{if } y_i \leq b_0 \end{cases}$$

$$\frac{\partial L}{\partial b_0} = \sum_{i=1}^n \begin{cases} -1 & \text{if } y_i > b_0 \\ 1 & \text{if } y_i \leq b_0 \end{cases}$$

Answer: at the median.



Let's try it in Python...



Summary

- **Models** can be written as $\hat{y} = f(x; \theta)$
- A **(training) loss function** measures how good or bad a fit is

$$\text{i.e. } L(\theta) = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

- **Recall 3 steps for supervised learning:**
 - 1) Choose a function form
 - 2) Choose a loss function
 - 3) Find the function form (i.e. parameters) which minimizes said loss