# Step 2: Goodness of Function

$y = b + w \cdot x_{cp}$

A set of function

Model
$f_1, f_2 \cdots$

Goodness of function f

Training Data

Loss function $L$:

Input: a function, output: how bad it is

Estimation error

$$L(f) = \sum_{n=1}^{10} \overline{\left( \hat{y}^n - \underline{f(x_{cp}^n)} \right)^2}$$

Sum over examples

Estimated y based on input function

$$L(w, b) = \sum_{n=1}^{10} \left( \hat{y}^n - \left( b + w \cdot x_{cp}^n \right) \right)^2$$
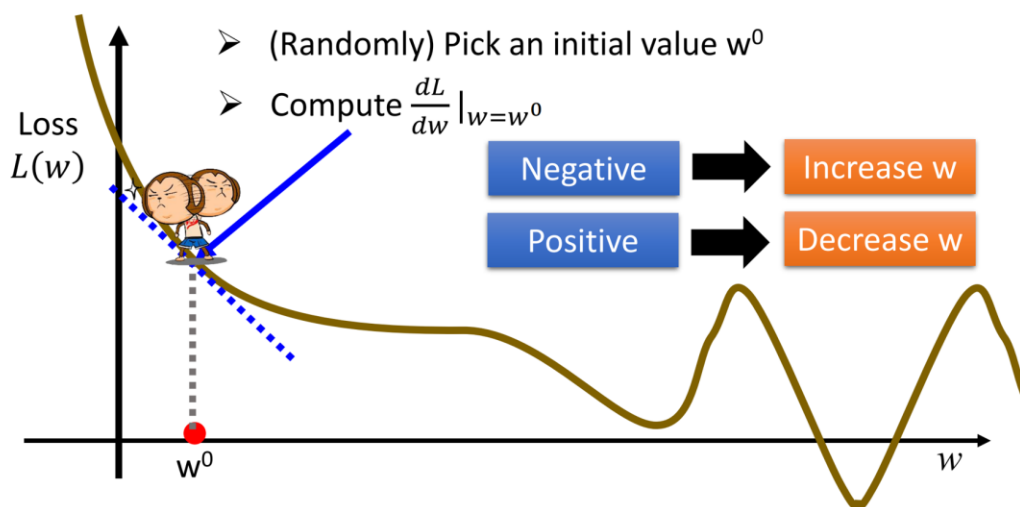
# Step 3: Gradient Descent

$$w^* = arg \min_w L(w)$$

- Consider loss function $L(w)$ with one parameter w:

  ➤ (Randomly) Pick an initial value $w^0$

  ➤ Compute $\frac{dL}{dw}\big|_{w=w^0}$

Loss $L(w)$

| Negative | ➡ | Increase w |
| Positive | ➡ | Decrease w |

$w^0$

$w$

# Step 3: Gradient Descent

$$w^* = arg \min_w L(w)$$

- Consider loss function $L(w)$ with one parameter w:

  ➢ (Randomly) Pick an initial value $w^0$

  ➢ Compute $\frac{dL}{dw}|_{w=w^0}$   $\quad w^1 \leftarrow w^0 - \eta \frac{dL}{dw}|_{w=w^0}$

  ➢ Compute $\frac{dL}{dw}|_{w=w^1}$   $\quad w^2 \leftarrow w^1 - \eta \frac{dL}{dw}|_{w=w^1}$

  …… Many iteration

Loss $L(w)$

Local minima

global minima

$w^0$   $w^1$  $w^2$   $w^T$   $w$

---

# Step 3: Gradient Descent

Loss

Very slow at the **plateau**

Stuck at saddle point

Stuck at local minima

$\partial L / \partial w \approx 0$

$\partial L / \partial w = 0$

$\partial L / \partial w = 0$

$L$

$w_1$

$w_2$

The value of the parameter w

# Back to step 2: Regularization

$$y = b + \sum w_i x_i$$

$$L = \sum_n \left( \hat{y}^n - \left( b + \sum w_i x_i \right) \right)^2 \quad \boxed{+\lambda \sum (w_i)^2}$$

➤ Smaller $w_i$ means … smoother

$$y = b + \sum w_i x_i$$

$$y + \sum w_i \Delta x_i = b + \sum w_i (x_i + \Delta x_i)$$

➤ We believe smoother function is more likely to be correct

Do you have to apply regularization on bias?