

DSC 40A

Theoretical Foundations of Data Science I

Last Time: Empirical Risk Minimization

- ▶ To learn, pick a **loss function** L and minimize the **empirical risk**:

$$R(h) = \frac{1}{n} \sum_{i=1}^n L(h, y_i)$$

- ▶ Absolute loss: $L_{\text{abs}}(h, y) = |h - y|$ (gives the **median**)
- ▶ Square loss: $L_{\text{sq}}(h, y) = (h - y)^2$ (gives the **mean**)
- ▶ **Key Point:** Tradeoffs to each loss function.

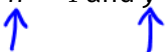
In This Video

We'll design our own loss function. We'll find that it's hard to minimize using the methods we've learned so far, which will motivate a new approach to minimizing functions.

Recommended Reading

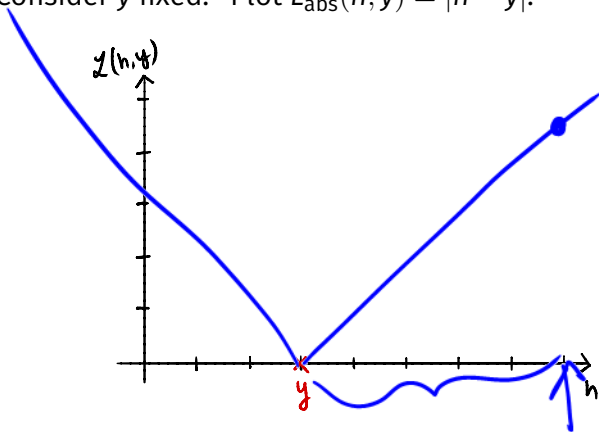
Course Notes: Chapter 1, Section 2

Loss Functions

- ▶ A loss function $L(h, y)$ quantifies how “bad” a prediction is.
- ▶ Example: take $h = 4$ and $y = 6$.

- ▶ Absolute loss: $L_{\text{abs}}(h, y) = |4 - 6| = 2$
- ▶ Square loss: $L_{\text{sq}}(h, y) = (4 - 6)^2 = 4$

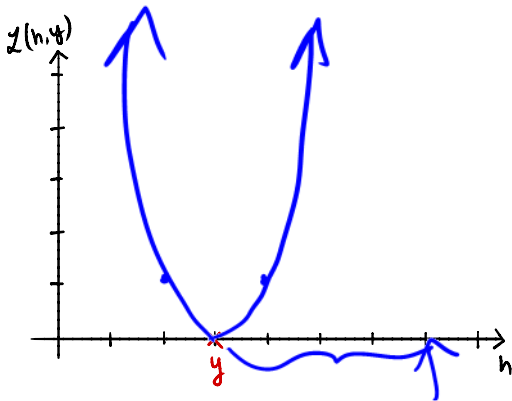
Plotting a Loss Function

- ▶ The plot of a loss function tells us how it treats outliers.
- ▶ Consider y fixed. Plot $L_{\text{abs}}(h, y) = |h - y|$:



Plotting a Loss Function

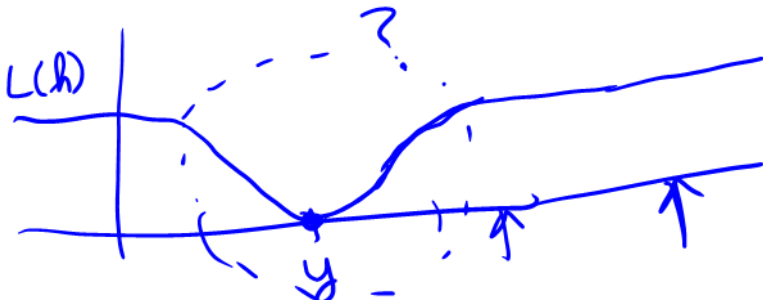
- ▶ The plot of a loss function tells us how it treats outliers.
- ▶ Consider y fixed. Plot $L_{\text{sq}}(h, y) = (h - y)^2$:



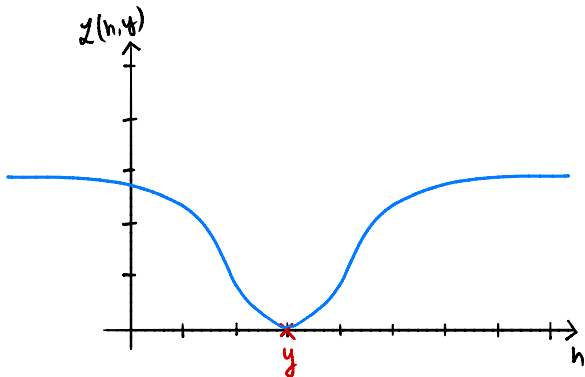
Question

Suppose L considers all outliers to be equally as bad. What would it look like far away from y ?

- a) flat
- b) rapidly decreasing
- c) rapidly increasing



A very insensitive loss




- ▶ We'll call this loss L_{UCSD} because it doesn't have a name.

L_{UCSD}

Question

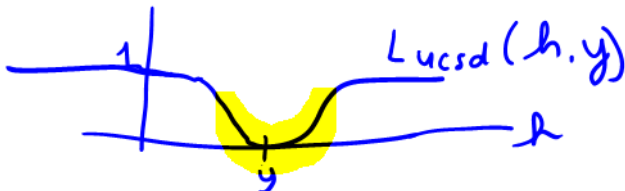
Which of these could be $L_{UCSD}(h, y)$?

~~a)~~ $e^{-\overset{\text{big}}{(h-y)^2}}$ 

b) $1 - e^{-\overset{\text{close to 0}}{(h-y)^2}}$

~~c)~~ $1 - (h-y)^2$

d) $1 - e^{-\overset{\text{big}}{(h-y)}}$
 close to 0

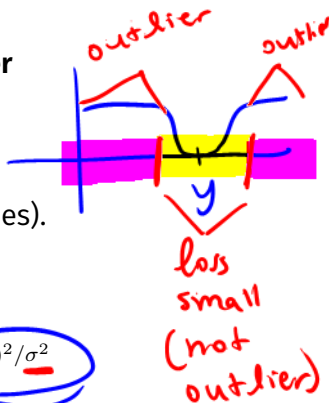
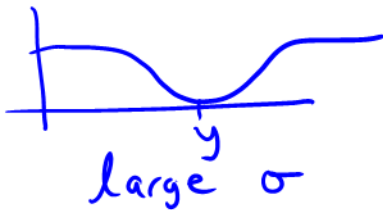
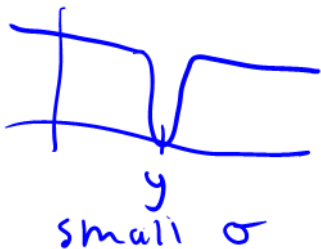


Adding a scale parameter

- ▶ Problem: L_{ucsd} has a fixed scale.
- ▶ Won't work for all data sets (e.g., salaries).
- ▶ Fix: add a **scale parameter**, σ :

$$L_{\text{ucsd}}(h, y) = 1 - e^{-(h-y)^2/\sigma^2}$$

similar form
as bell curve



Empirical Risk Minimization

- ▶ We have salaries y_1, \dots, y_n .
- ▶ To find prediction, ERM says to minimize the mean loss:

$$\begin{aligned} R_{\text{ucsd}}(h) &= \frac{1}{n} \sum_{i=1}^n L_{\text{ucsd}}(h, y_i) \\ &= \frac{1}{n} \sum_{i=1}^n \left[1 - e^{-(h-y_i)^2/\sigma^2} \right] \end{aligned}$$

Let's plot R_{ucsd}

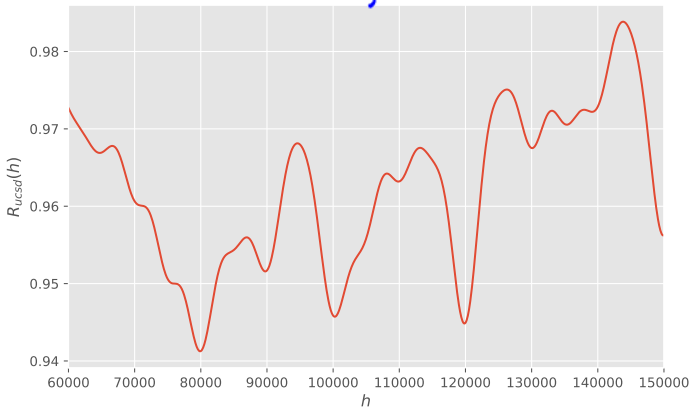
- ▶ Recall:

$$R_{\text{ucsd}}(h) = \frac{1}{n} \sum_{i=1}^n \left[1 - e^{-(h-y_i)^2/\sigma^2} \right]$$

- ▶ Once we have data y_1, \dots, y_n and a scale σ , we can plot $R_{\text{ucsd}}(h)$
- ▶ We'll use full StackOverflow data ($n = 1121$)
- ▶ Let's try several scales, σ .

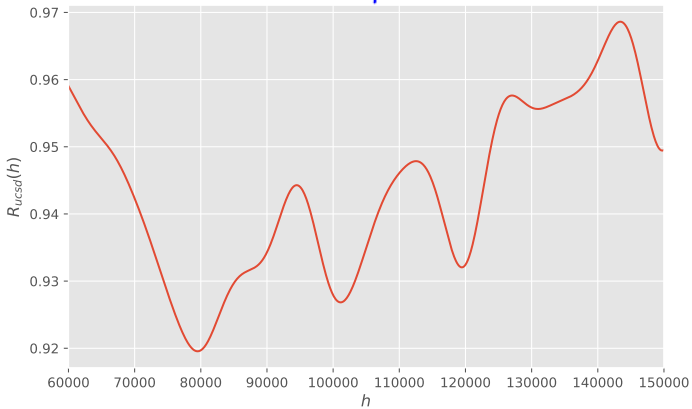
Plot of R_{ucsd}

$\sigma = 3000$



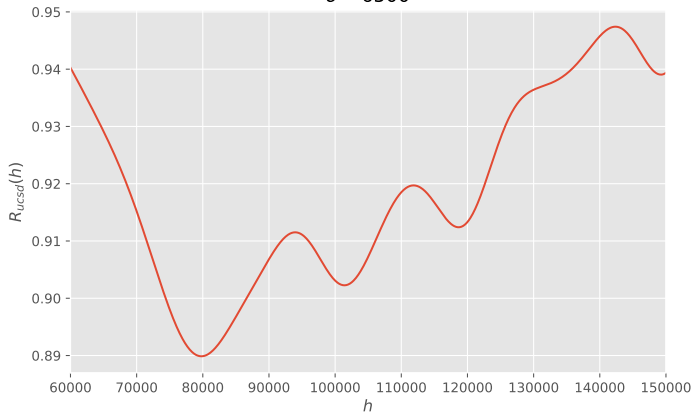
Plot of R_{ucsd}

$\sigma = 4500$



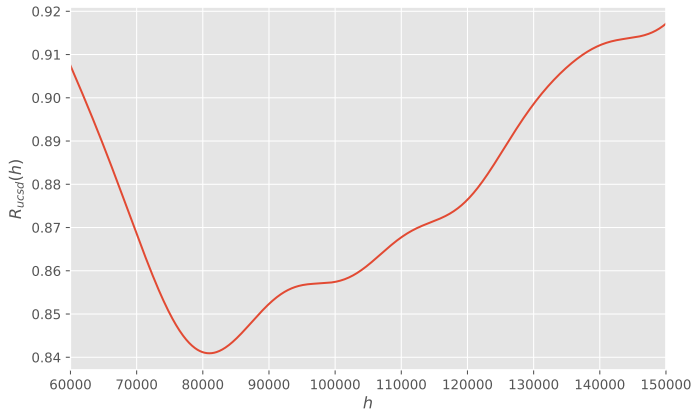
Plot of R_{ucsd}

$\sigma = 6500$



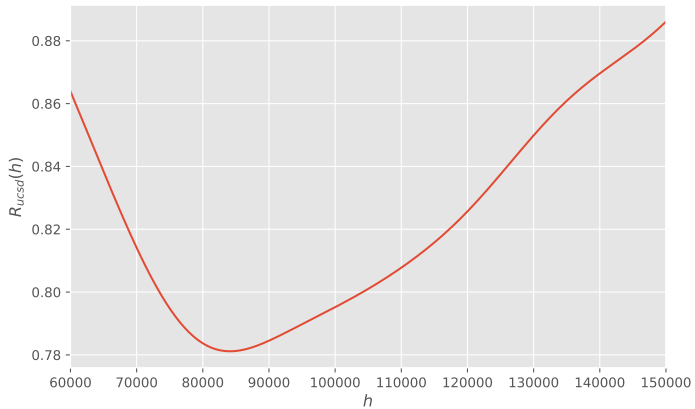
Plot of R_{ucsd}

$\sigma = 10000$



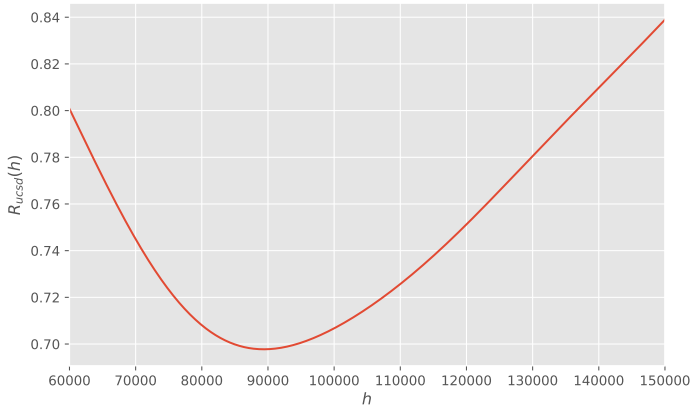
Plot of R_{ucsd}

$\sigma = 14500$



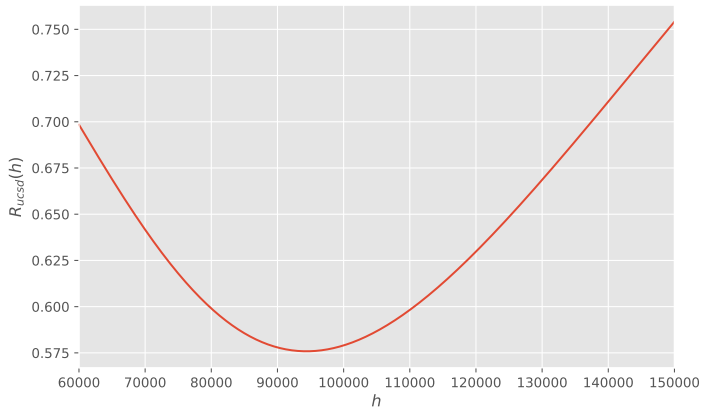
Plot of R_{ucsd}

$\sigma = 21000$



Plot of R_{ucsd}

$\sigma = 32000$



Minimizing R_{ucsd}

- ▶ To make prediction, we find h^* minimizing $R_{\text{ucsd}}(h)$.
- ▶ R_{ucsd} is **differentiable**.
- ▶ To minimize: take derivative, set to zero, solve.

Step 1) Taking the derivative

$$\begin{aligned}\frac{dR_{UCSD}}{dh} &= \frac{d}{dh} \left(\frac{1}{n} \sum_{i=1}^n \left[1 - e^{-(h-y_i)^2/\sigma^2} \right] \right) \\ &= \frac{1}{n} \sum_{i=1}^n \frac{d}{dh}(1) - \frac{d}{dh} \left(e^{-(h-y_i)^2/\sigma^2} \right) \\ &= \frac{1}{n} \sum_{i=1}^n - e^{-\boxed{(h-y_i)^2/\sigma^2}} \cdot -2 \boxed{(h-y_i)}/\sigma^2 \cdot 1 \\ &= \frac{2}{n\sigma^2} \sum_{i=1}^n e^{-(h-y_i)^2/\sigma^2} \cdot (h-y_i)\end{aligned}$$

Step 2) Setting to zero and solving

- ▶ We found:

$$\frac{dR_{\text{ucsd}}}{dh}(h) = \frac{2}{n\sigma^2} \sum_{i=1}^n (h - y_i) \cdot e^{-(h-y_i)^2/\sigma^2}$$

- ▶ Now we just set to zero and solve for h :

$$0 = \frac{2}{n\sigma^2} \sum_{i=1}^n (h - y_i) \cdot e^{-(h-y_i)^2/\sigma^2}$$

- ▶ We **can** calculate derivative, but we **can't** solve for h ; we're stuck again.

Summary

- ▶ We created our own loss function, which was designed to treat all outliers in much the same way.
- ▶ Our loss function was differentiable, but we still couldn't minimize it
- ▶ **Next Time:** We'll invent a general algorithm called **gradient descent** for minimizing differentiable functions.