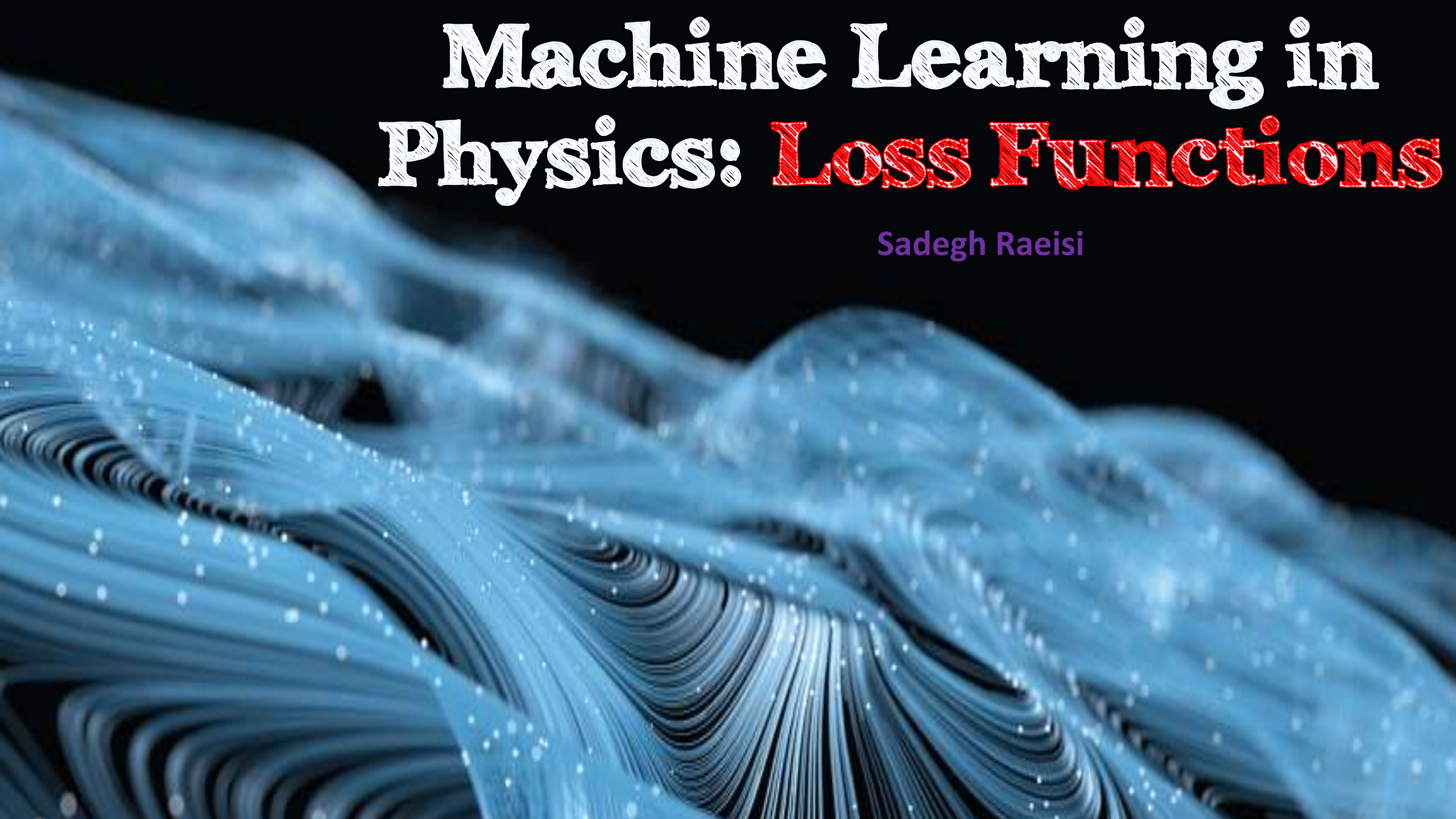
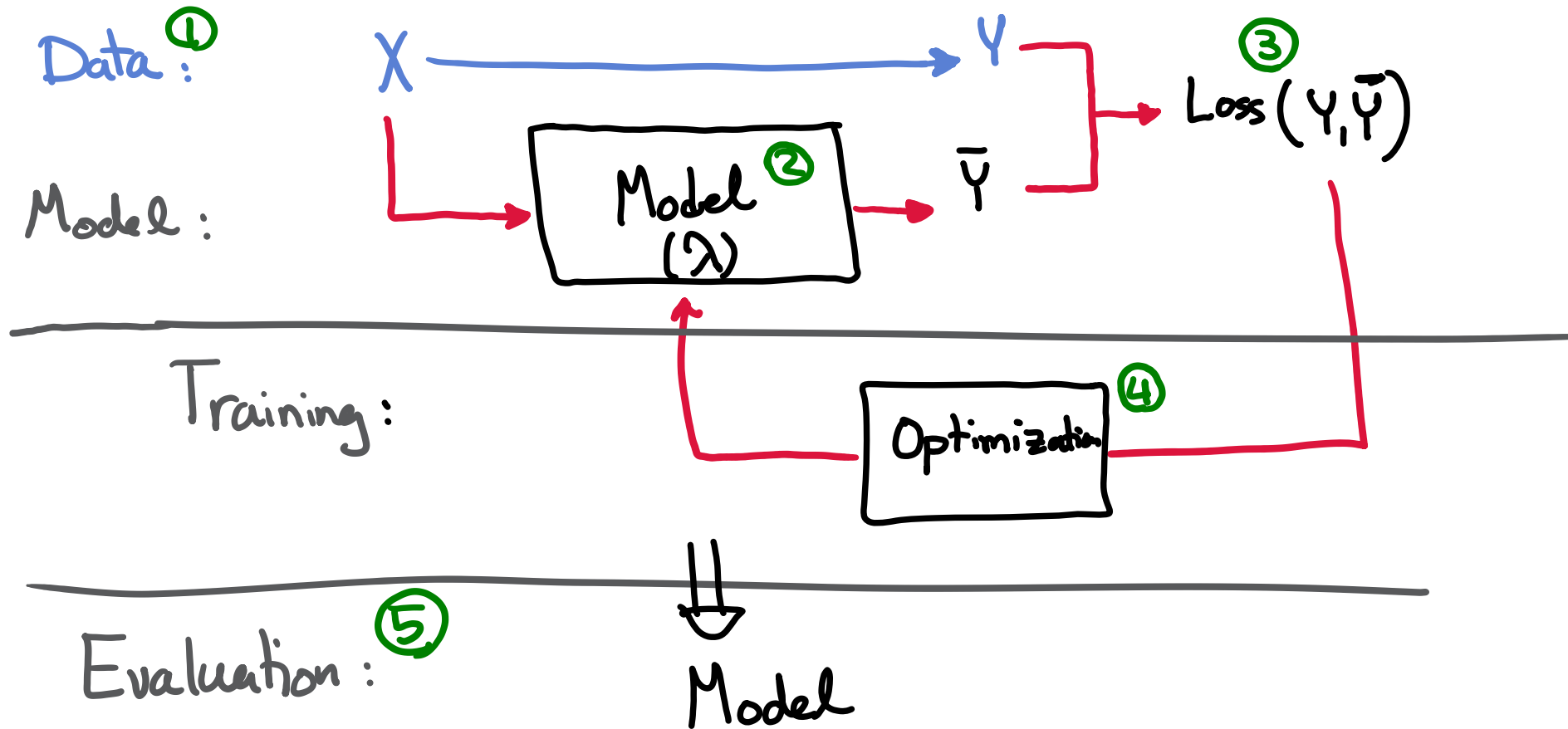


Machine Learning in Physics: **Loss Functions**

Sadegh Raeisi



Supervised: Ingredients



Outline



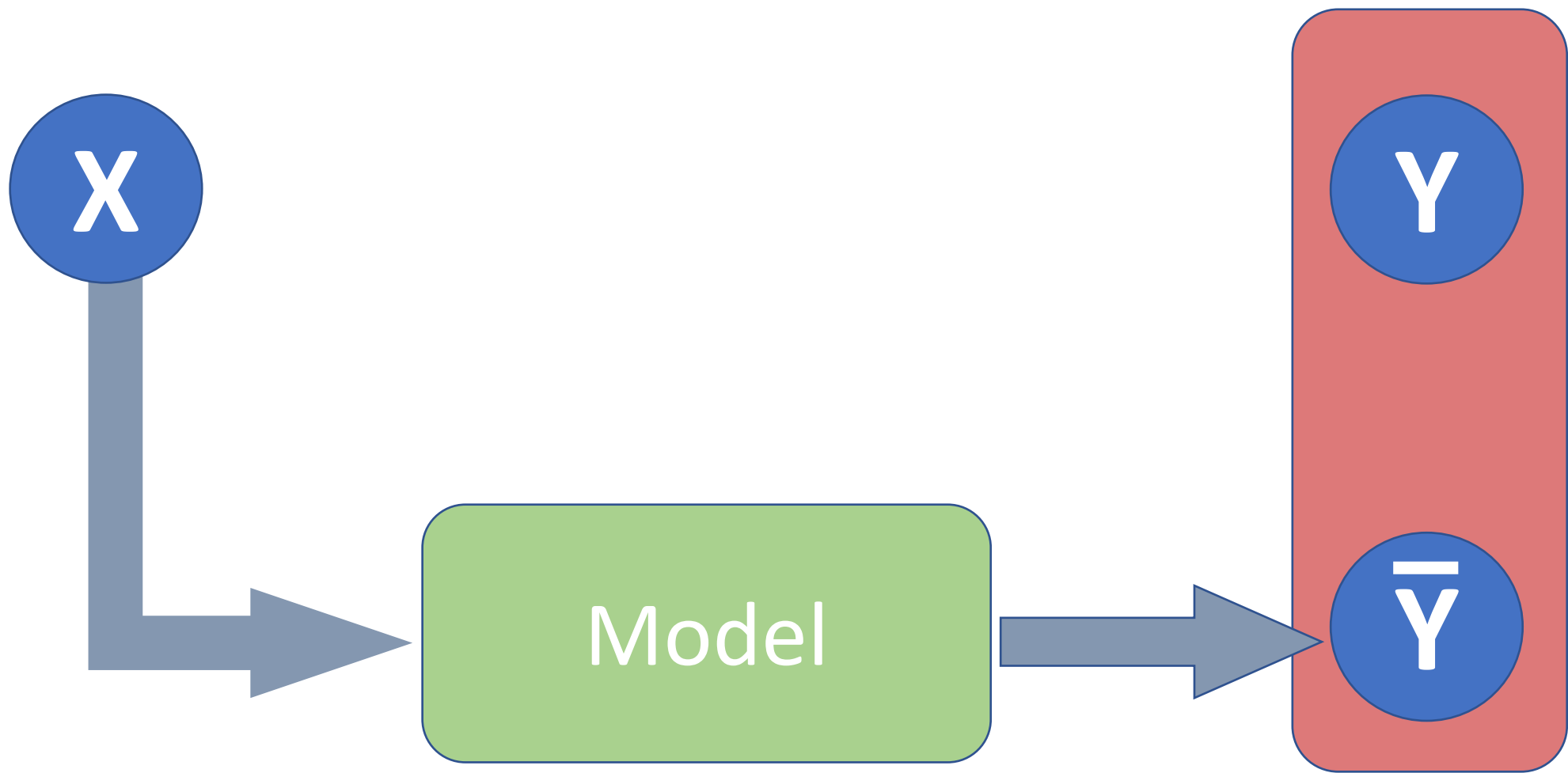
Concept

Regression

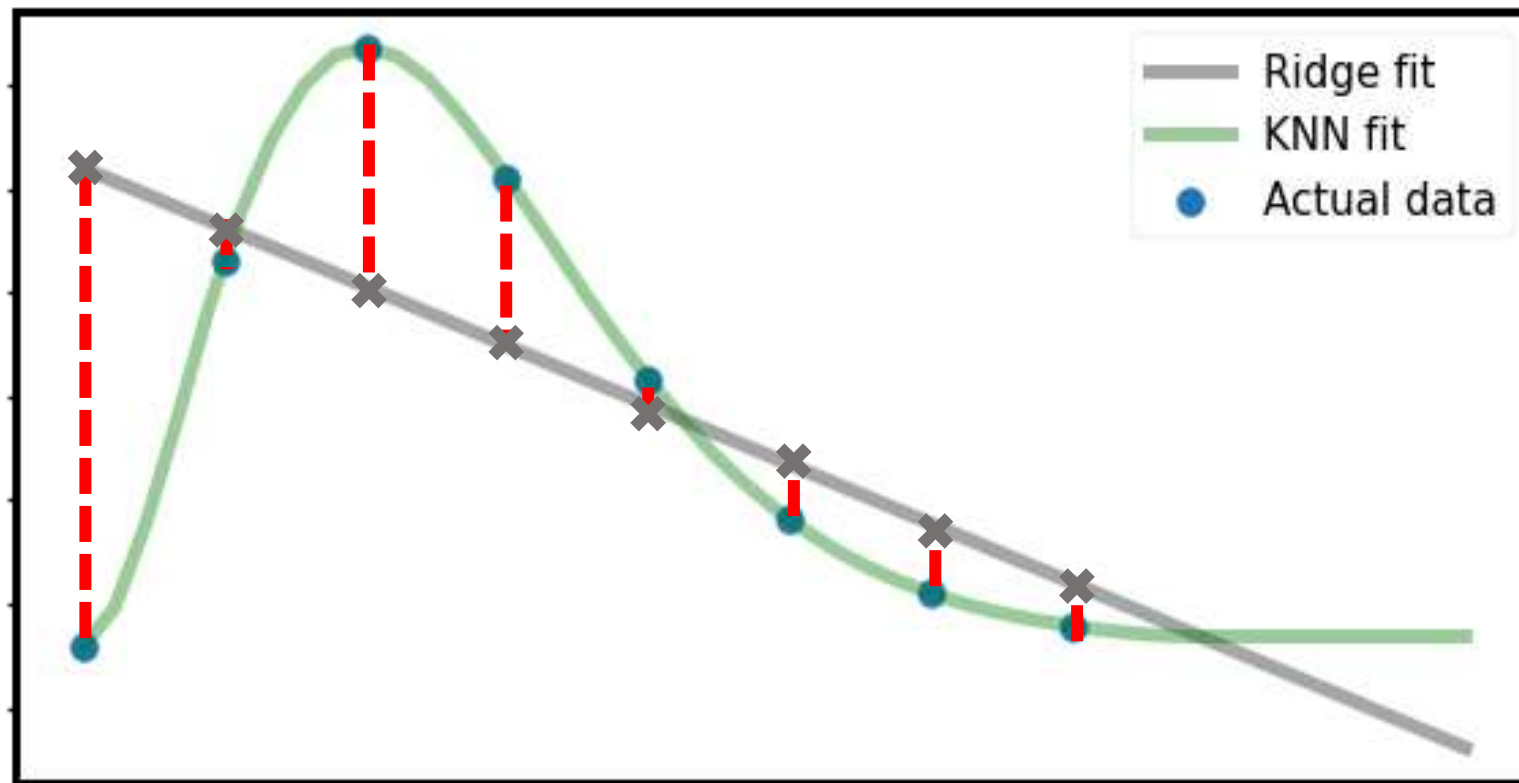
Classification

Other Loss functions

Concept



How close are the predictions?

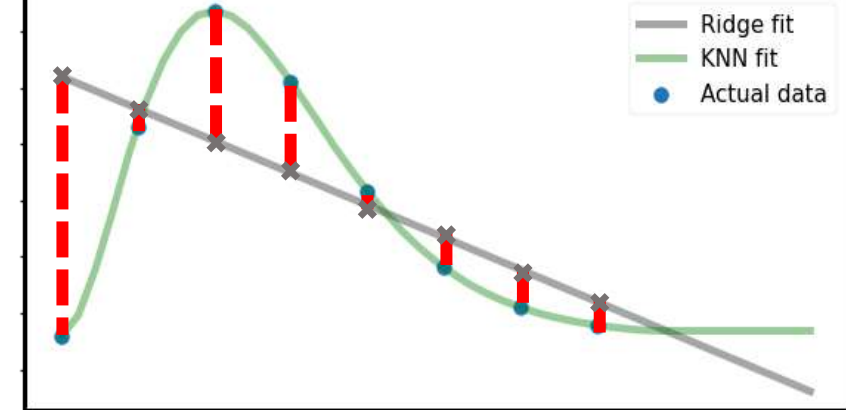


How can we quantify the difference?

$$\text{Dist}(\textcircled{Y}, \textcircled{\bar{Y}})$$

Regression

Minkowski distance



$$D(Y, \bar{Y}) = \left(\sum_i |Y^i - \bar{Y}^i|^p \right)^{\frac{1}{p}}$$

Minkowski distance

$$D_1(Y, \bar{Y}) = \sum_i |Y^i - \bar{Y}^i|$$

$$D_2(Y, \bar{Y}) = \sqrt{\sum_i |Y^i - \bar{Y}^i|^2}$$

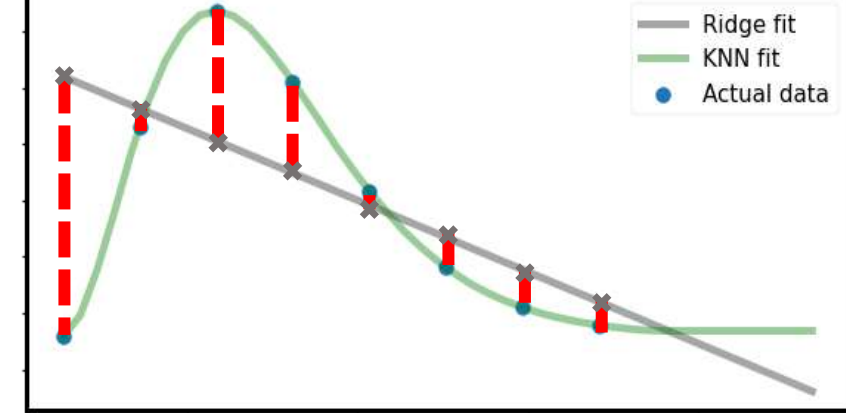
$$D_\infty(Y, \bar{Y}) = \max_i |Y^i - \bar{Y}^i|$$



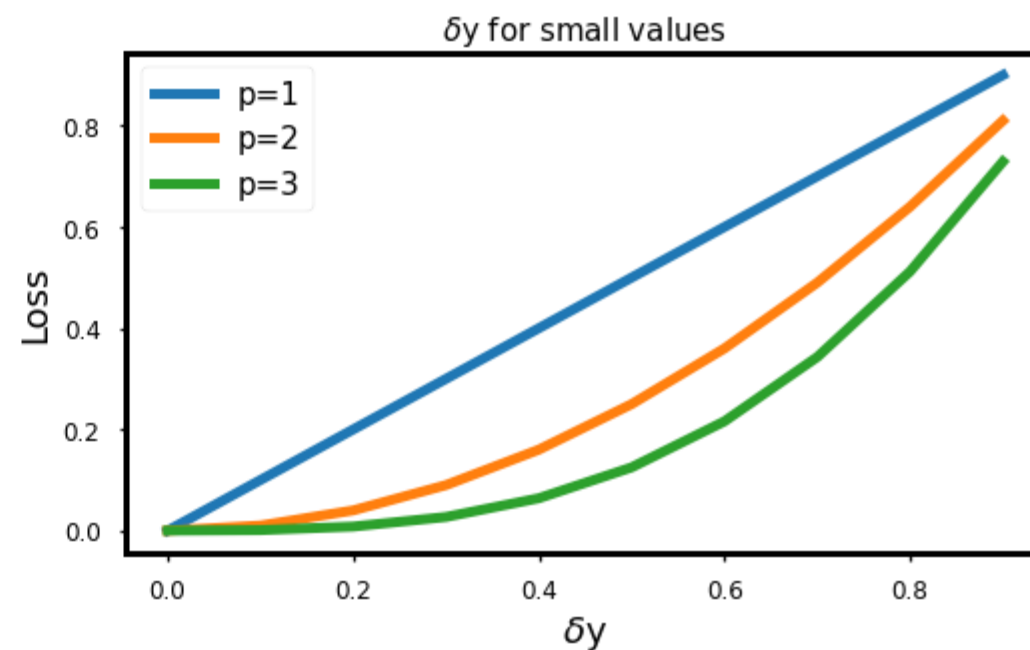
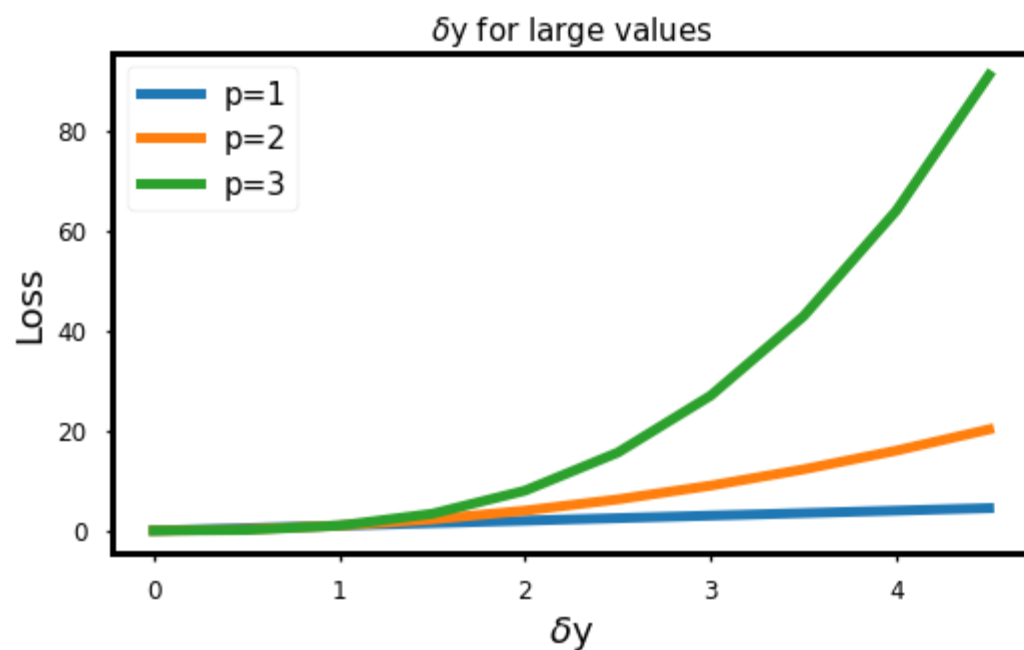
Penalizes larger deviations more



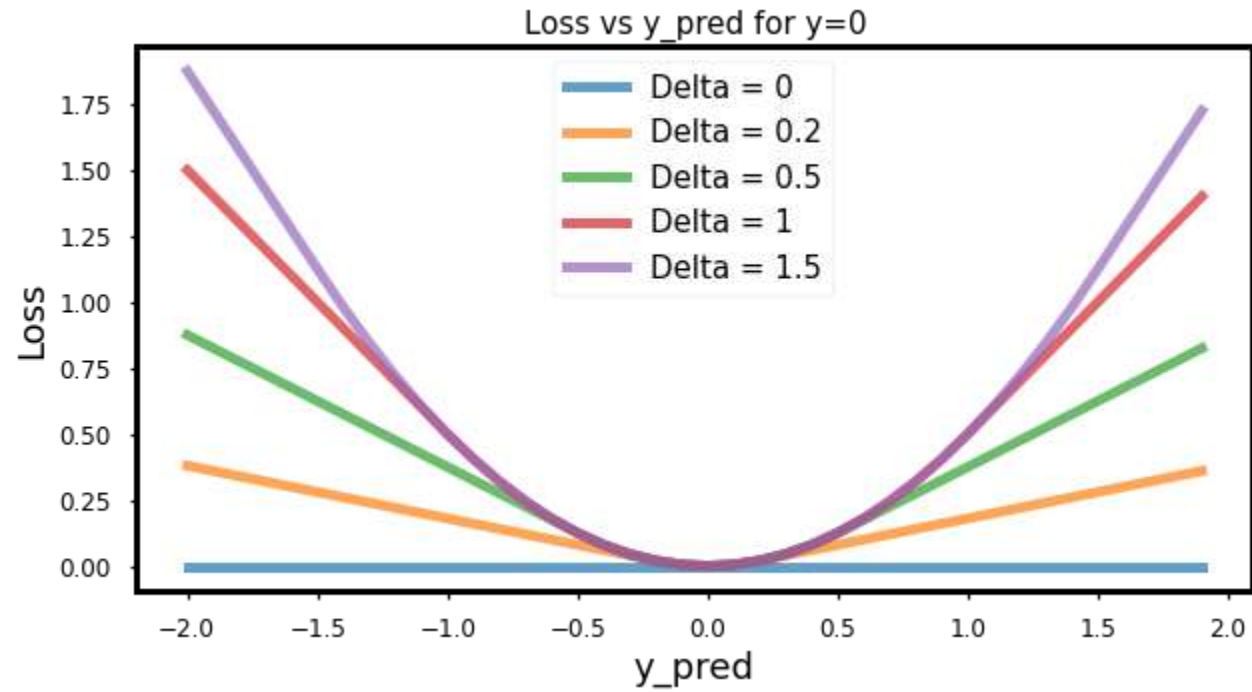
Worst case



Minkowski distance



Huber: the best of both worlds

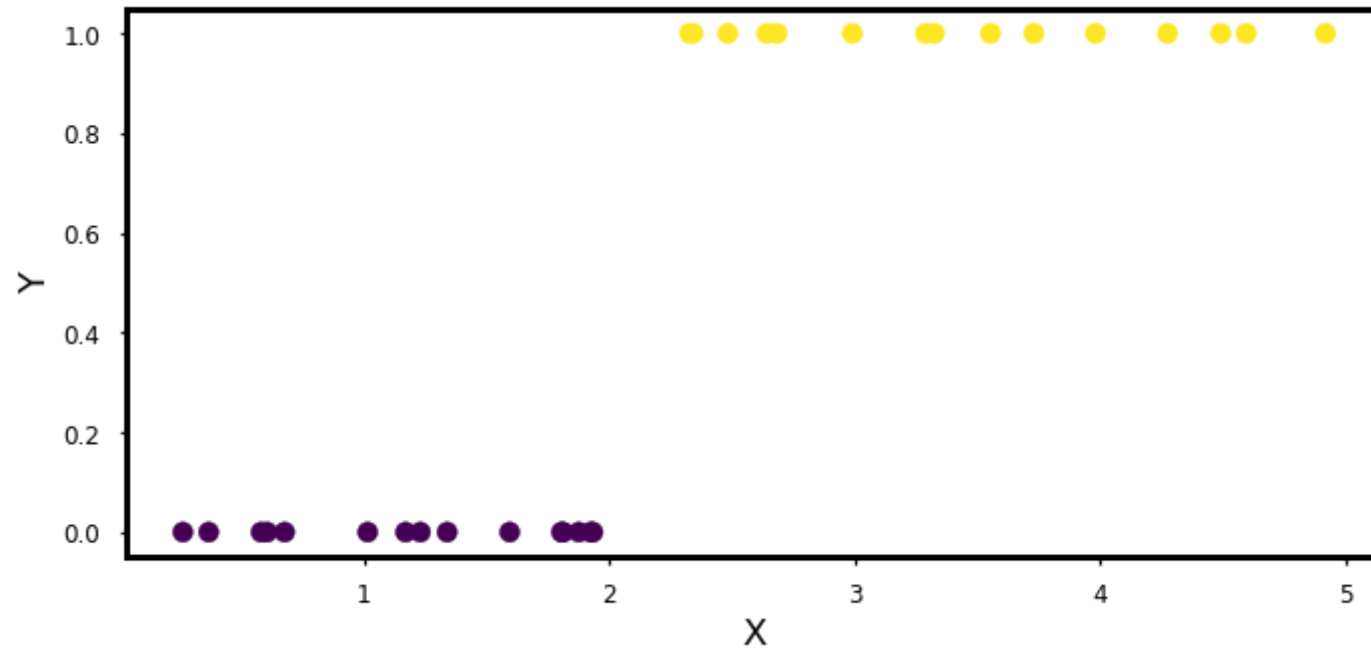


What other functions?

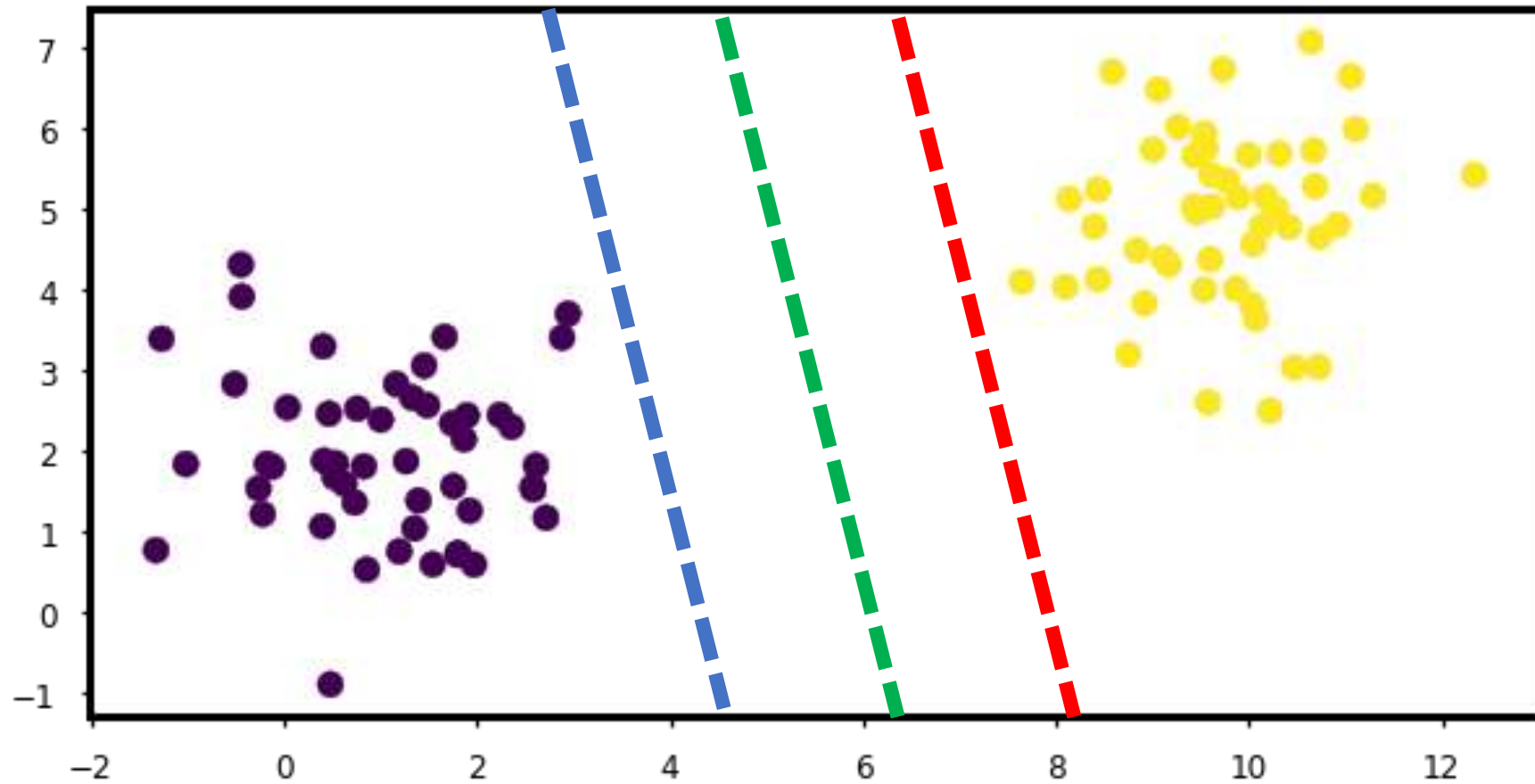
Take a look at [scikit-learn](#)!

Classification

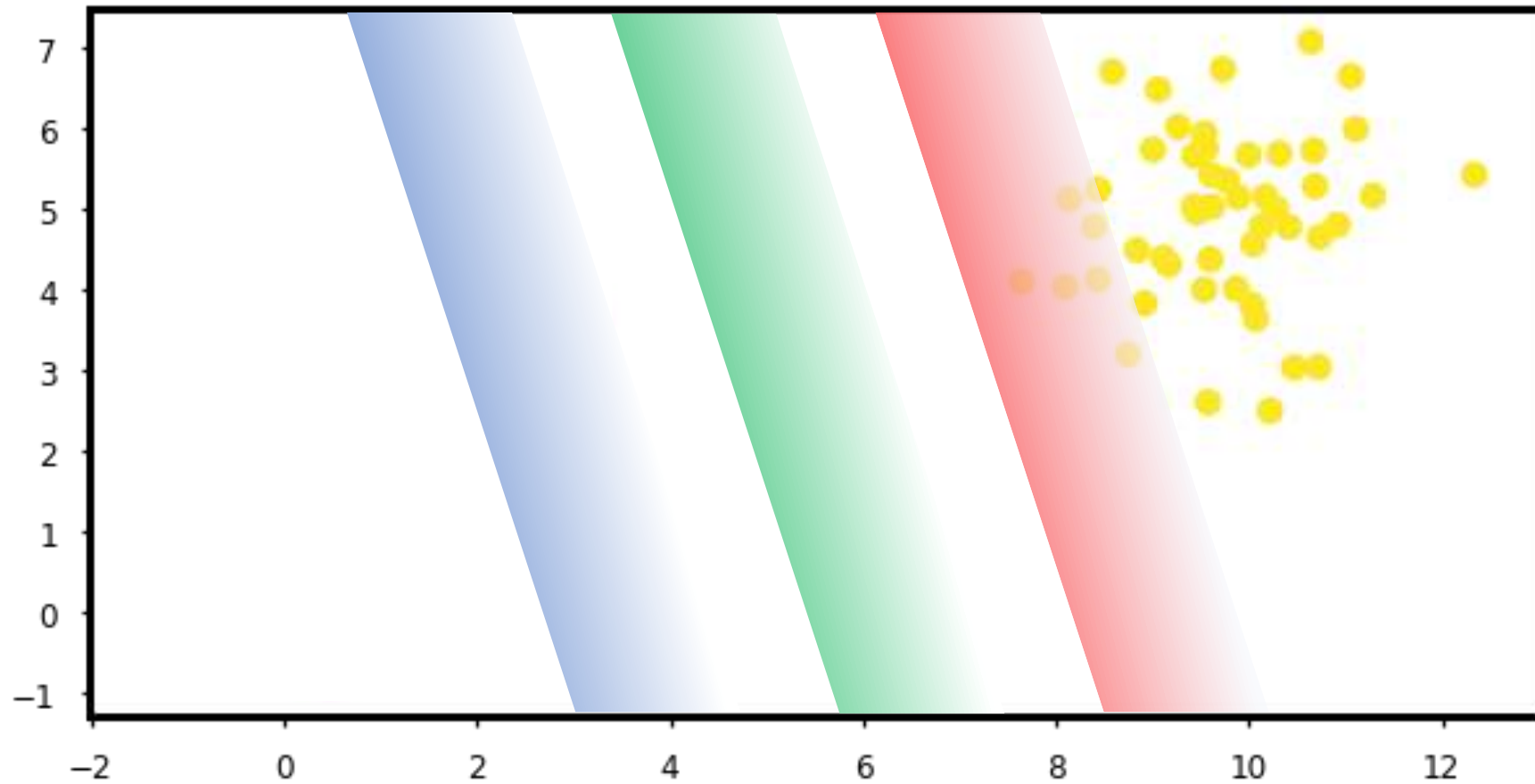
We can use regression loss functions



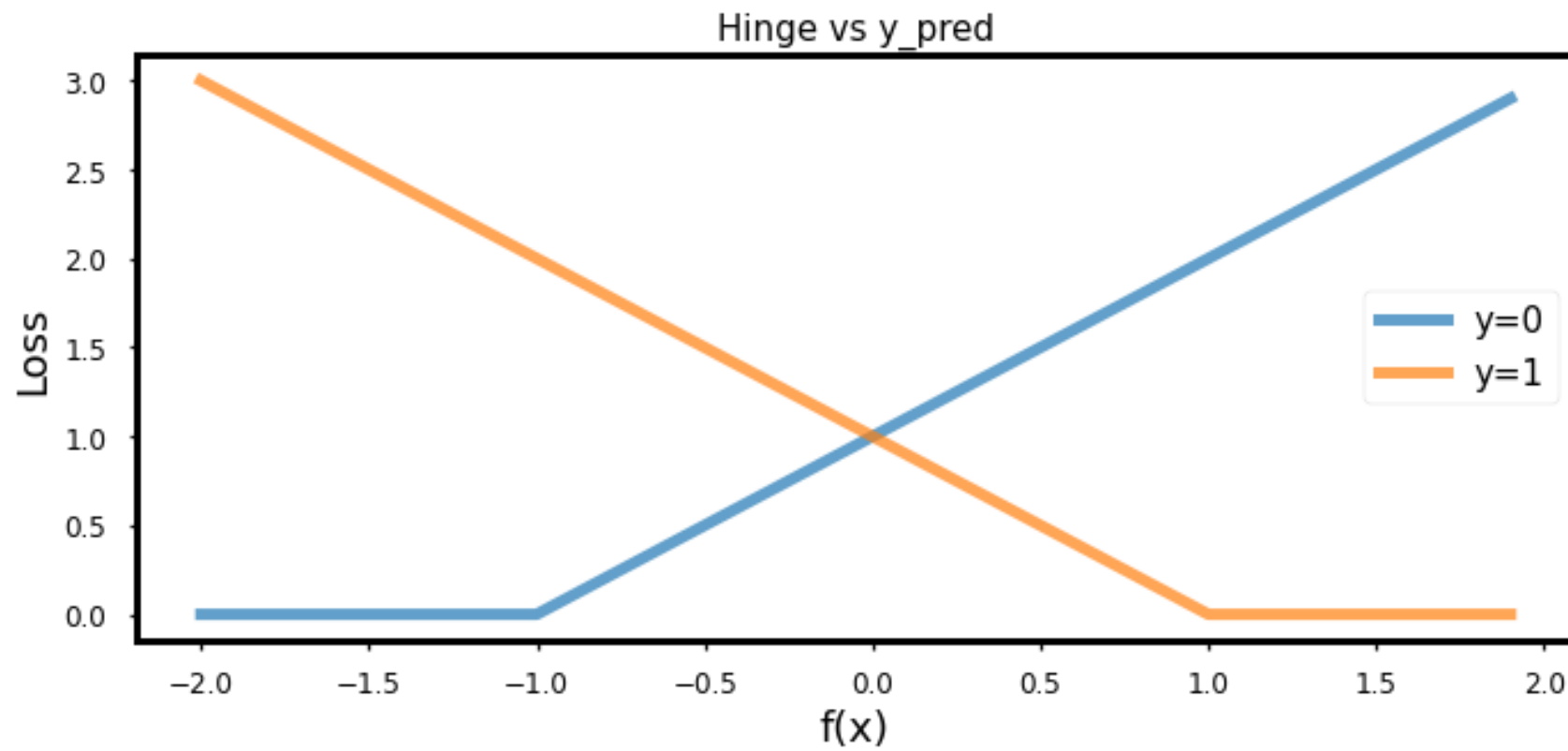
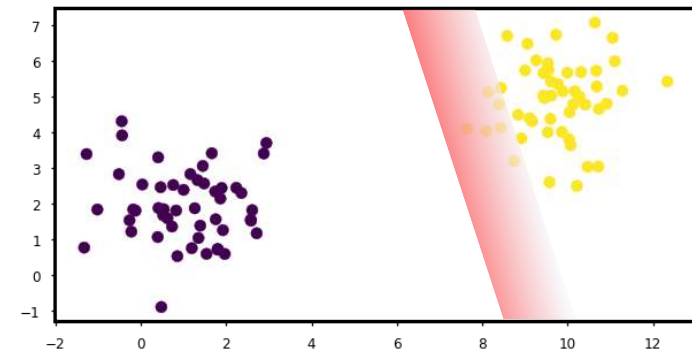
Hinge



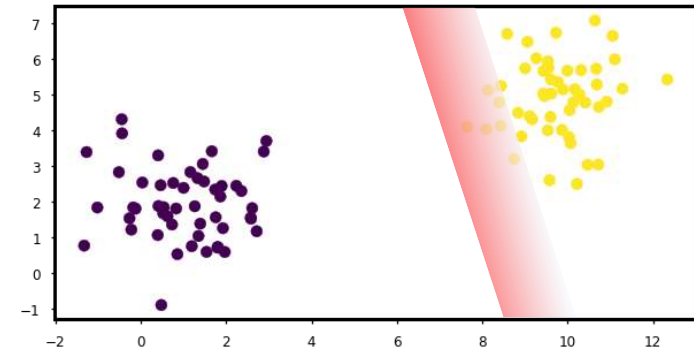
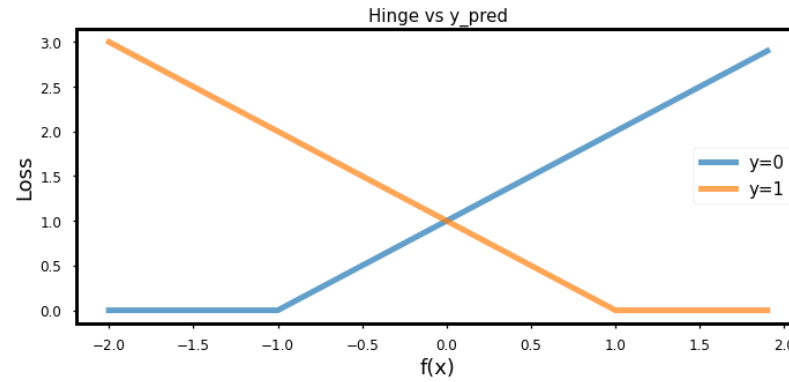
Hinge



Hinge



Hinge

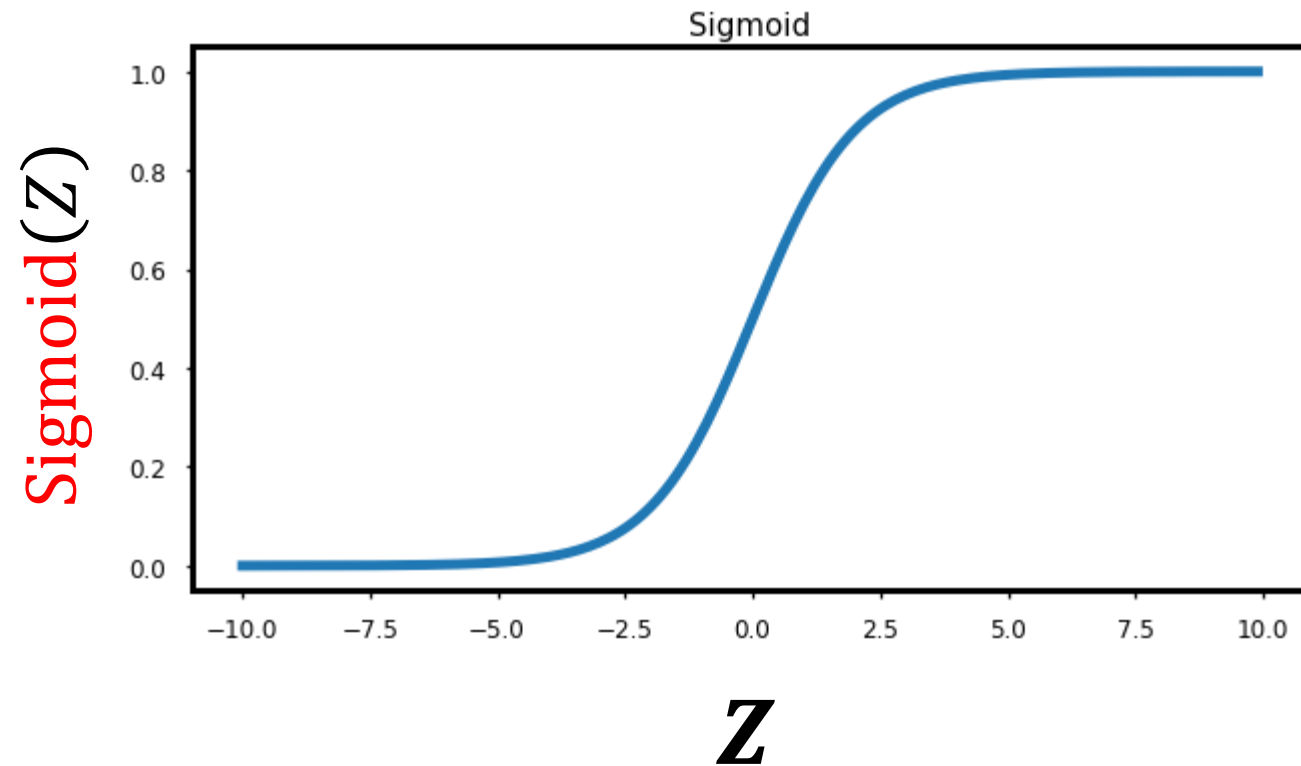


$$D_h(Y, \bar{Y}) = \max(0, 1 - Y^i * f_w(X^i))$$

Note that f_w is the decision boundary, not the class.
Also note that $Y^i = \pm 1$ and not $\{0,1\}$.

Probability:

Logistic Regression

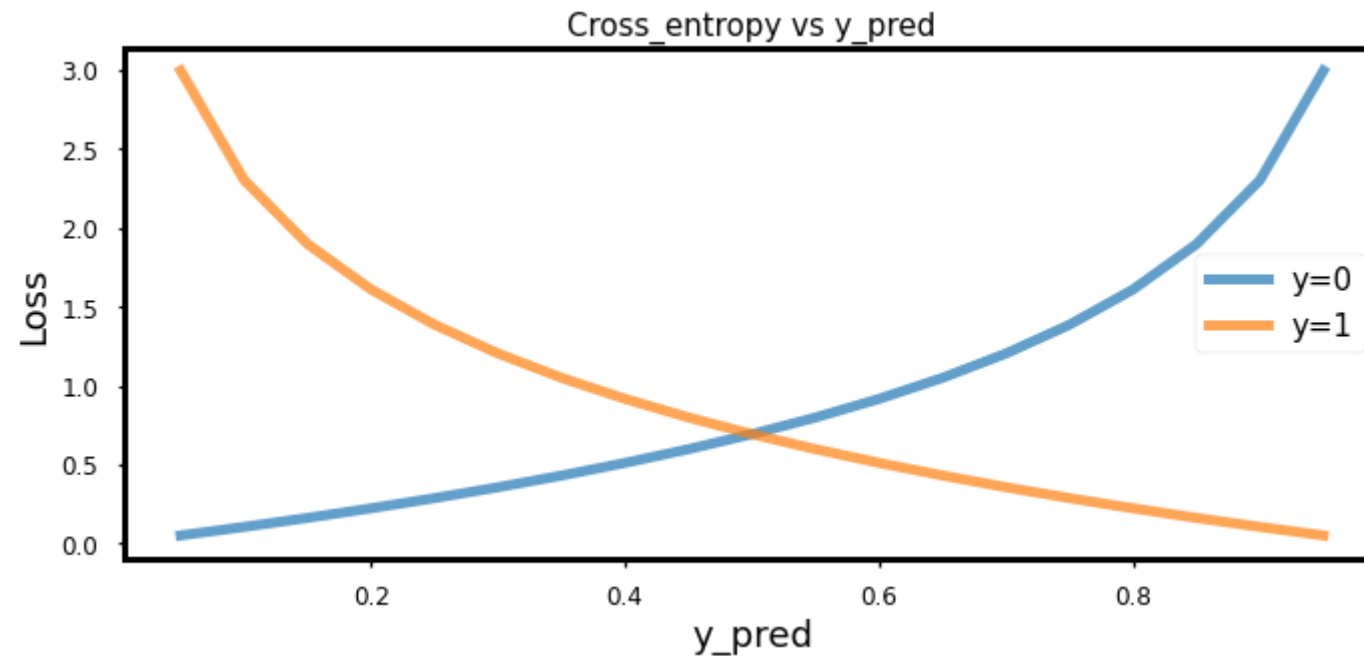


Probability: Cross-entropy

$$D_{CE}(Y, \bar{Y}) = - \sum_i Y^i \log \bar{Y}^i + (1 - Y^i) \log(1 - \bar{Y}^i)$$

Probability: Cross-entropy

$$D_{CE}(Y, \bar{Y}) = - \sum_i Y^i \log \bar{Y}^i + (1 - Y^i) \log(1 - \bar{Y}^i)$$



Other loss functions

What would be a good loss function
for a multi-class problem?

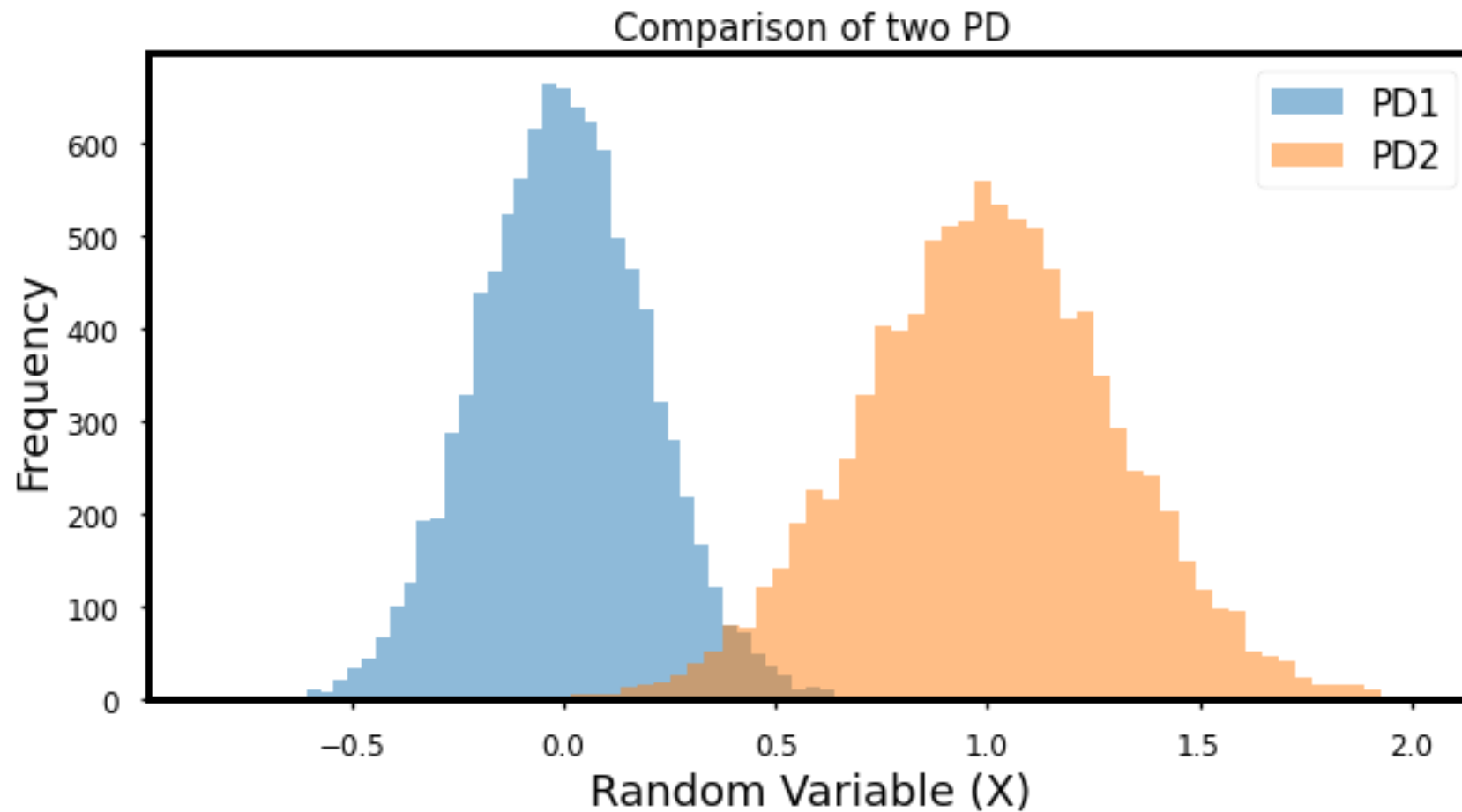
- $Y \in \{-1, 0, 1\}$

What would be a good
loss function for $Y \in \{0, 1\}^{\otimes n}$?

Example:

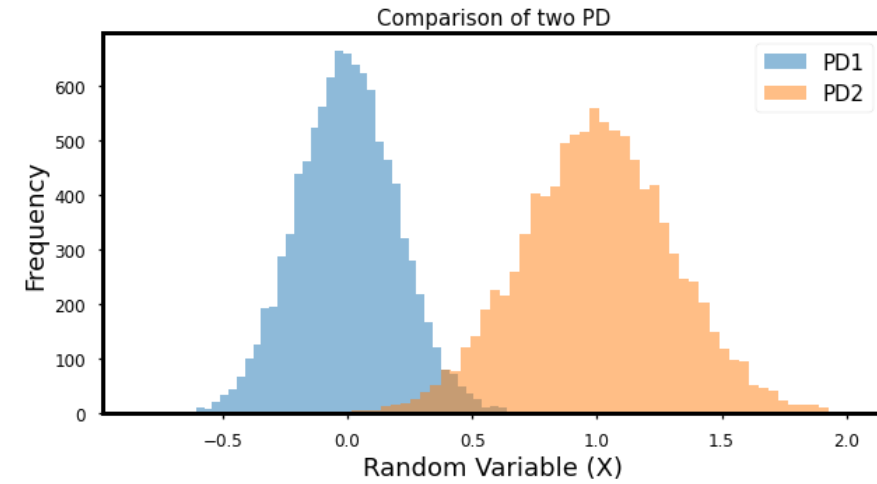
$$Y^1 = (0, 0, 1, 0, 1)$$

Probability Distributions



Probability Distributions

- Kullback–Leibler divergence (Relative Entropy)



$$D_{KL}(PD_1 || PD_2) = \sum_X PD_1(X) \log\left(\frac{PD_1(X)}{PD_2(X)}\right)$$

What would be a good
loss function for clustering?

**The right choice
for the
Loss function**

How do we choose the right loss function?

- Convexity
- Objective

How does this choice affect the model?

See the NB.

So far ...

