

Empirical Risk Minimization

Alexander Jung

Learning Goals:

- know about notion of **expected loss or risk**
- know that **average loss approximates risk**
- know about **empirical risk minimization**
- be aware of **design choices data/model/loss** and their **effect on ERM** methods

What is ML About ?

fit **models** to **data** to make
predictions or forecasts !

Data. Model. Loss.

data: set of datapoints (x,y)

model: set of hypothesis maps $h(\cdot)$

loss: quality measure $L((x,y),h)$

Data

	Year	m	d	Time	Time zone	Maximum temperature (degC)	Minimum temperature (degC)
0	2020	2	1	00:00	UTC	3.0	1.9
1	2020	2	2	00:00	UTC	4.9	2.4
2	2020	2	3	00:00	UTC	2.6	-0.4
3	2020	2	4	00:00	UTC	-0.2	-3.7
4	2020	2	5	00:00	UTC	2.5	-4.2
5	2020	2	6	00:00	UTC	2.4	-4.7
6	2020	2	7	00:00	UTC	1.2	-5.5
7	2020	2	8	00:00	UTC	2.7	0.2
8	2020	2	9	00:00	UTC	3.9	2.6

Data.

$$\mathcal{D} = \{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(m)}, y^{(m)})\}.$$

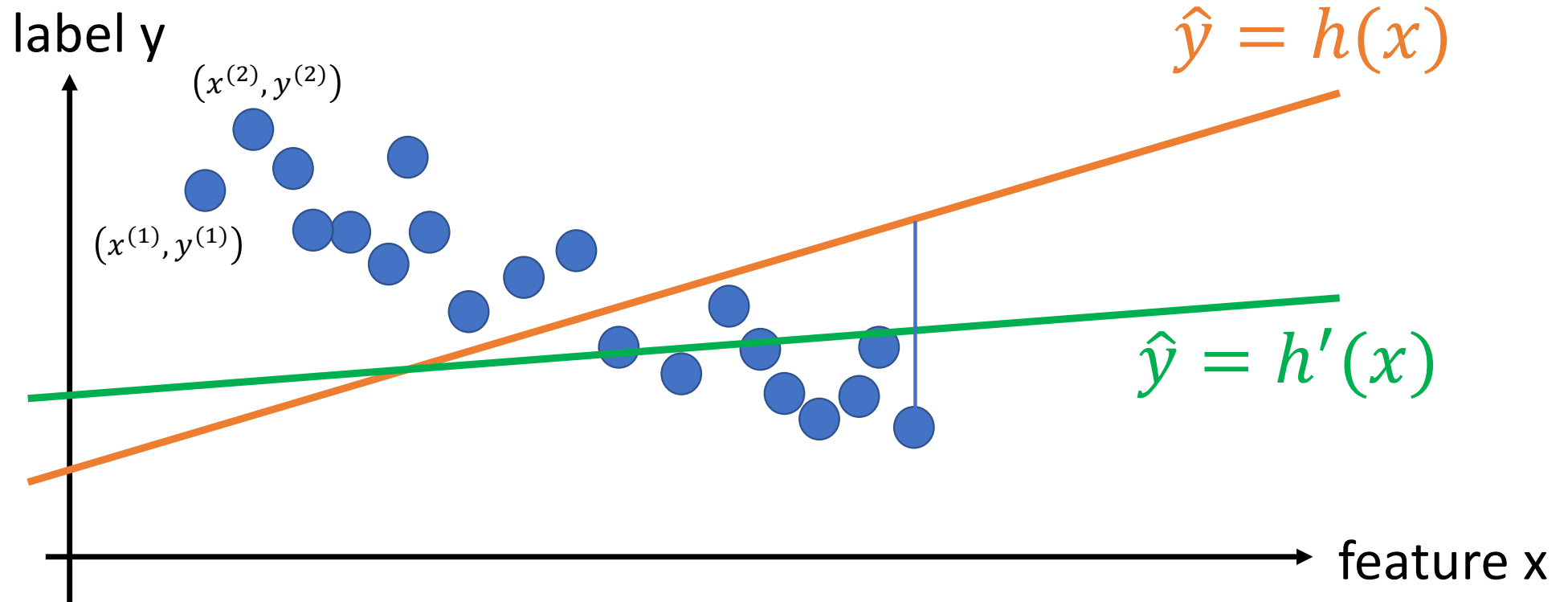
	Year	m	d	Time	Time zone	Maximum temperature (degC)	Minimum temperature (degC)
0	2020	2	1	00:00	UTC	3.0	1.9
1	2020	2	2	00:00	UTC	4.9	2.4
2	2020	2	3	00:00	UTC	2.6	-0.4
3	2020	2	4	00:00	UTC	-0.2	-3.7
4	2020	2	5	00:00	UTC	2.5	-4.2
5	2020	2	6	00:00	UTC	2.4	-4.7
6	2020	2	7	00:00	UTC	1.2	-5.5
7	2020	2	8	00:00	UTC	2.7	0.2
8	2020	2	9	00:00	UTC	3.9	2.6

stack feature vecs into matrix

$$\mathbf{X} = (\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)})^T \in \mathbb{R}^{m \times n}$$

stack labels into vector

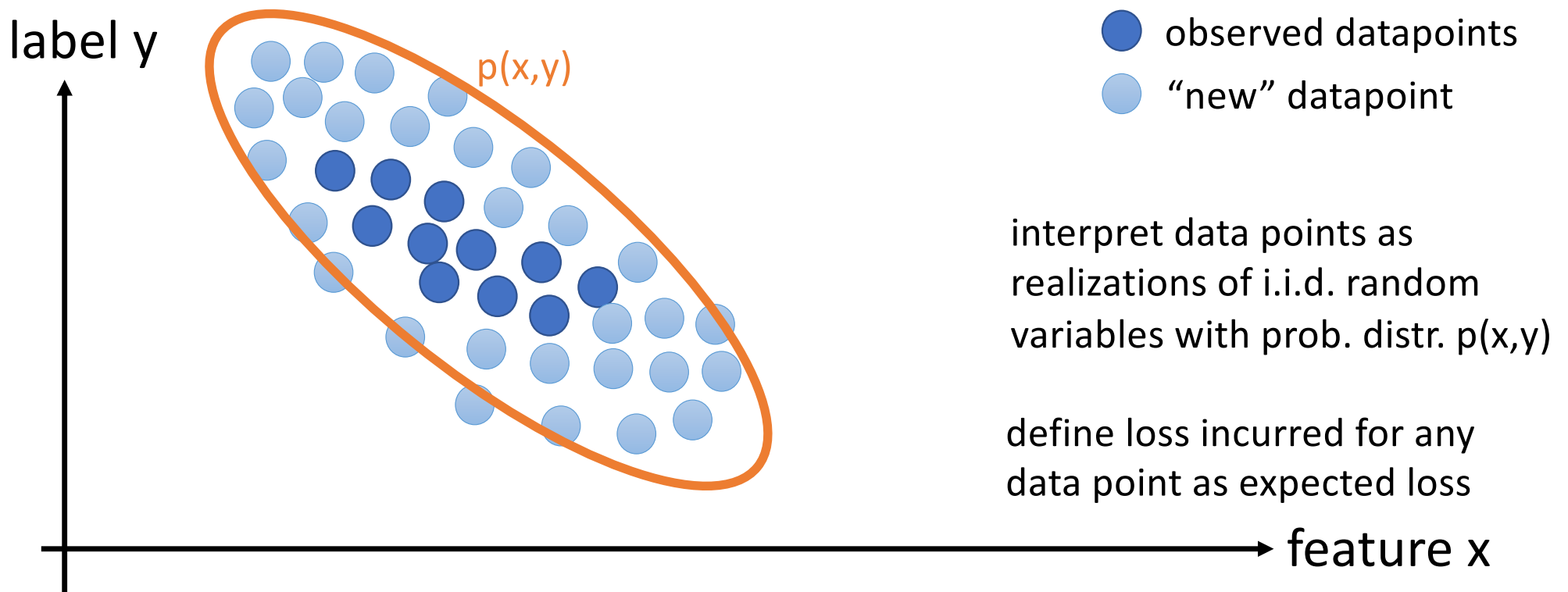
$$\mathbf{y} = (y^{(1)}, \dots, y^{(m)})^T \in \mathbb{R}^m$$



Machine Learning.

find hypothesis in model that incurs
smallest loss when predicting label of
any datapoint

What is Any Datapoint?



Expected Loss or Risk

$$\mathbb{E}\{L((\mathbf{x}, y), h)\} := \int_{\mathbf{x}, y} L((\mathbf{x}, y), h) dp(\mathbf{x}, y) \quad (2.14)$$

note: to compute this expectation
we need to know the probability distribution
 $p(\mathbf{x}, y)$ of datapoints (\mathbf{x}, y)

Empirical Risk

IDEA: approximate expected loss by average loss on some datapoints (training set)

$$\mathcal{D} = \{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(m)}, y^{(m)})\}.$$

$$\mathbb{E}\{L((\mathbf{x}, y), h)\} \approx (1/m) \sum_{i=1}^m L((\mathbf{x}^{(i)}, y^{(i)}), h) \text{ for sufficiently large sample size } m. \quad (2.17)$$

with the average loss or **empirical risk**

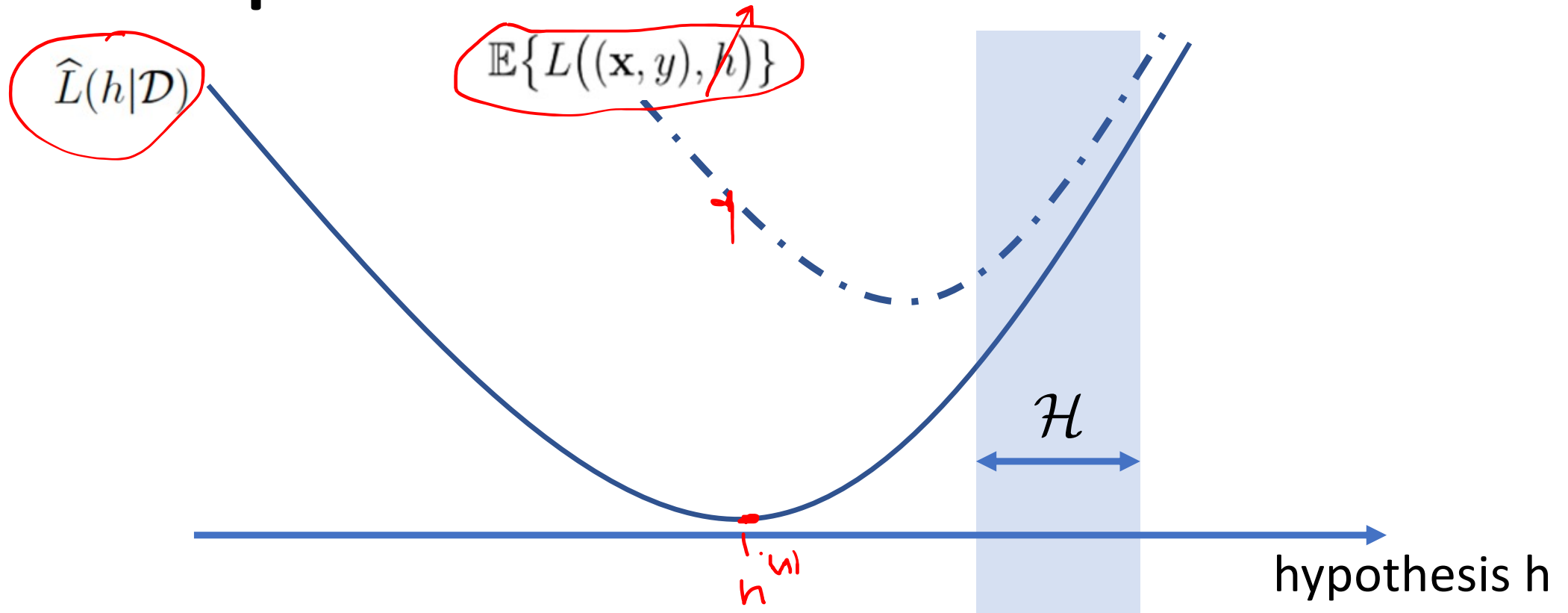
$$\hat{L}(h|\mathcal{D}) = (1/m) \sum_{i=1}^m L((\mathbf{x}^{(i)}, y^{(i)}), h). \quad (2.16)$$

Empirical Risk Minimization

$$\hat{h} \in \operatorname{argmin}_{h \in \mathcal{H}} \hat{L}(h|\mathcal{D})$$

$$\stackrel{(2.16)}{=} \operatorname{argmin}_{h \in \mathcal{H}} (1/m) \sum_{i=1}^m L((\mathbf{x}^{(i)}, y^{(i)}), h).$$

Empirical Risk Minimization



ERM for Parametrized Models

learnt (optimal) parameter vector

$$\hat{\mathbf{w}} = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} f(\mathbf{w})$$

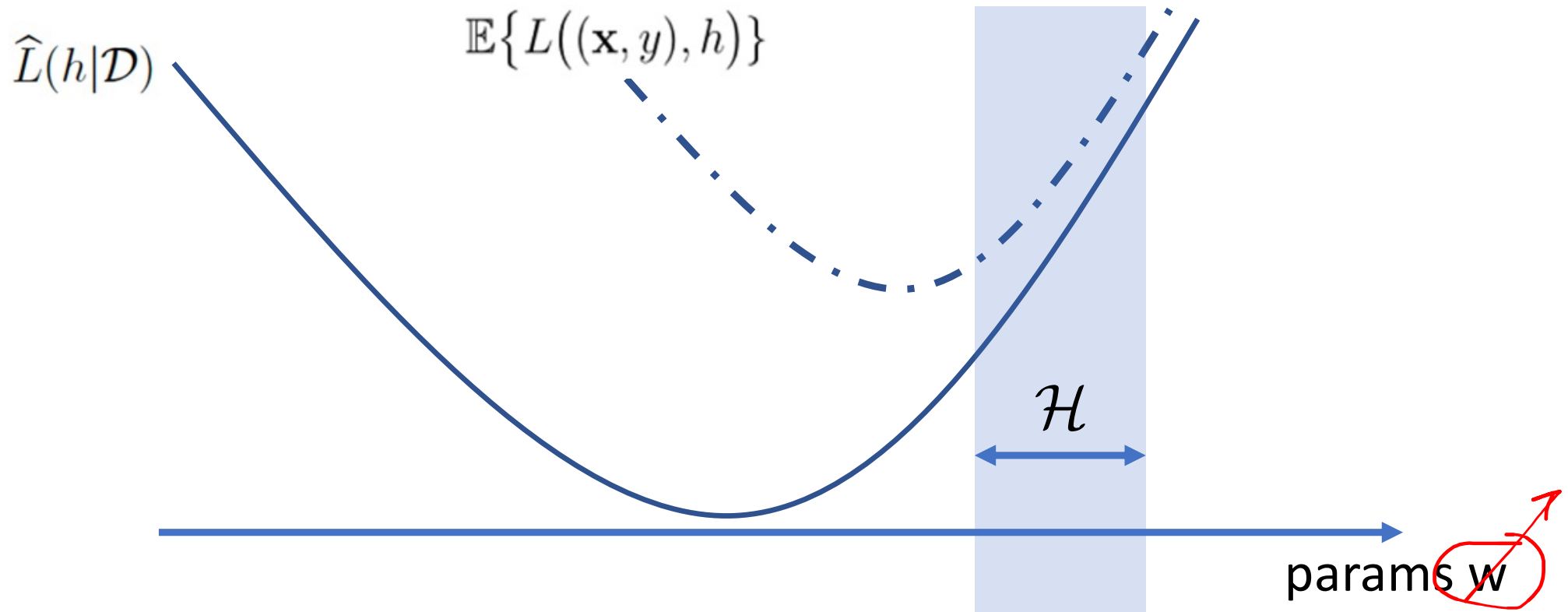
$$h(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$$

loss incurred by $h(\cdot)$
for i -th data point

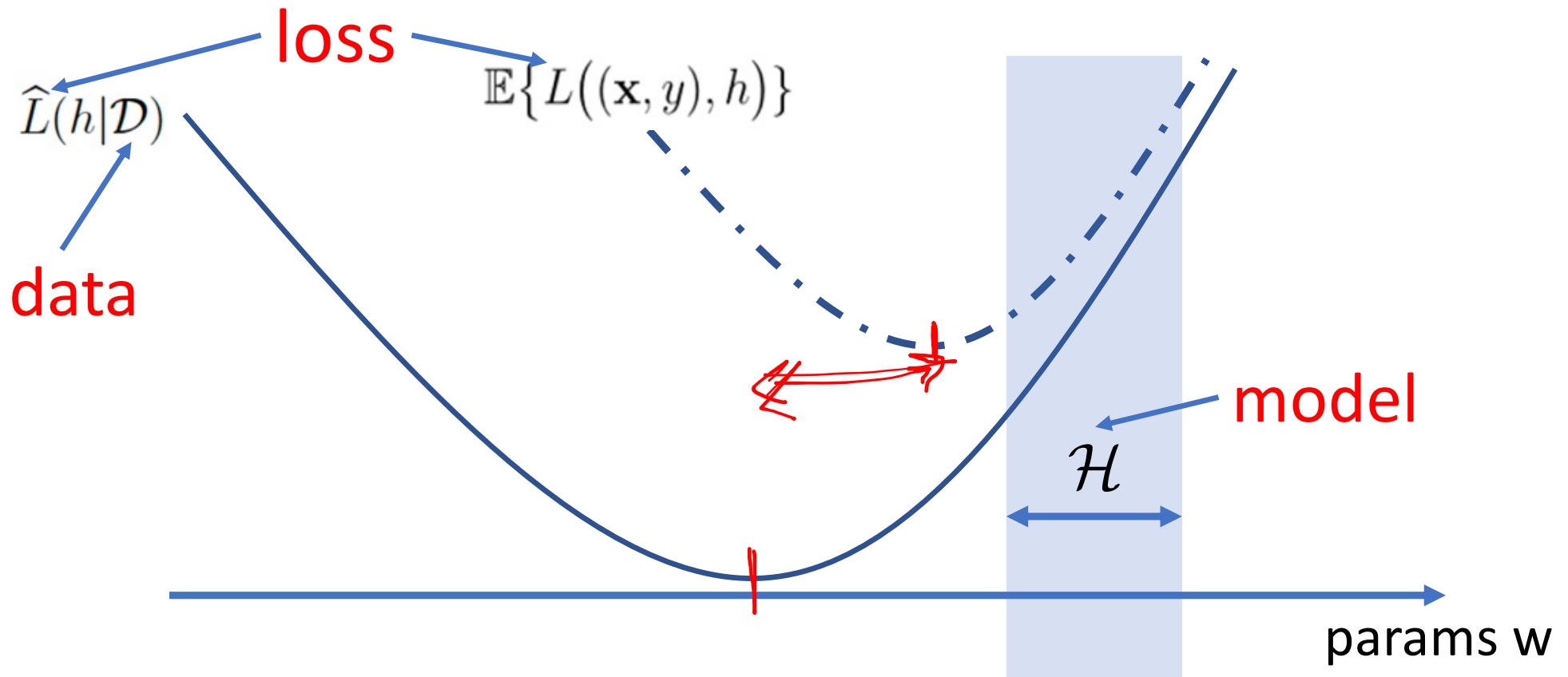
$$\text{with } f(\mathbf{w}) := \underbrace{\left(\frac{1}{m} \sum_{i=1}^m L((\mathbf{x}^{(i)}, y^{(i)}), h^{(\mathbf{w})}) \right)}_{\hat{L}(h^{(\mathbf{w})} | \mathcal{D})}$$

$\hat{L}(h^{(\mathbf{w})} | \mathcal{D})$ average loss or
empirical risk

ERM for Param. Models



Design Choices in ERM

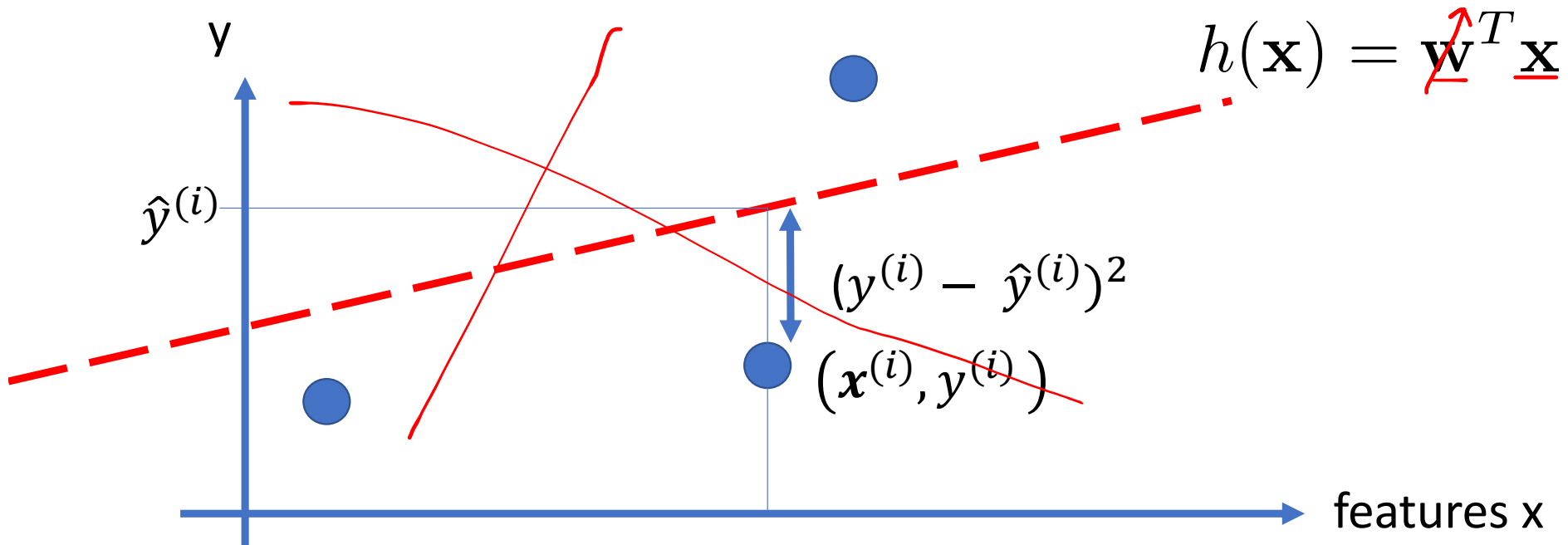


Design Choice: Model and Data

Linear Regression

- datapoints characterized by feature vector and numeric label
- model consists of linear hypothesis maps
- squared error loss

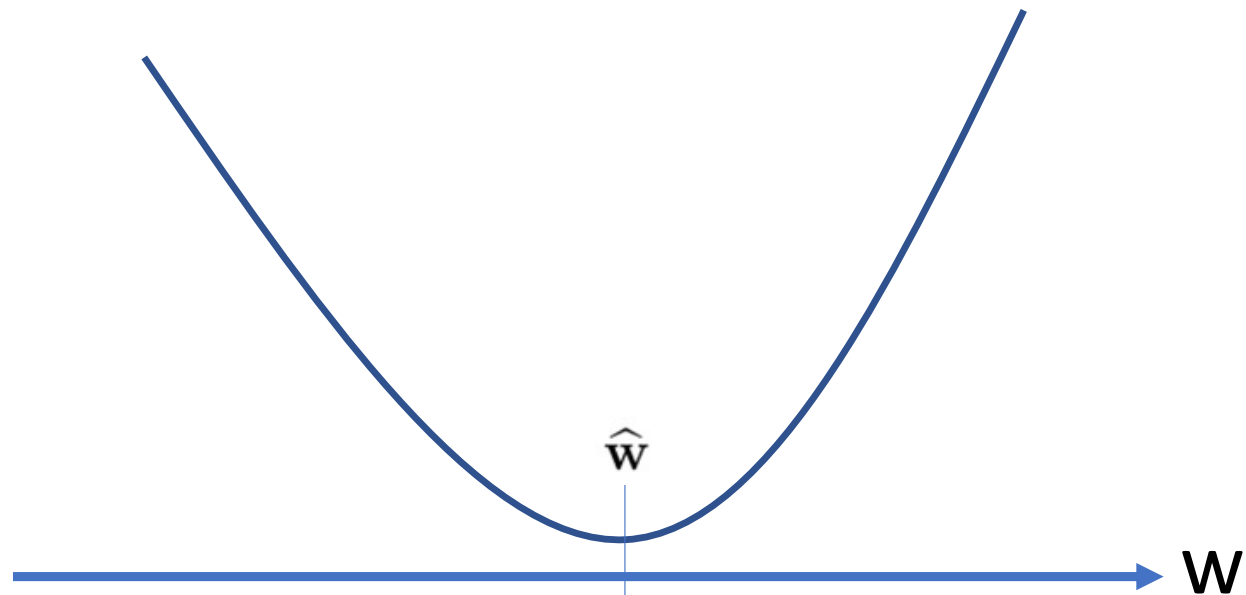
Linear Regression



choose parameter/weight vector \mathbf{w} to minimize average squared error loss

ERM for Linear Regression

$$\hat{\mathbf{w}} = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^n} (1/m) \sum_{m=1}^m (y^{(i)} - \mathbf{w}^T \mathbf{x}^{(i)})^2. \quad (4.5)$$



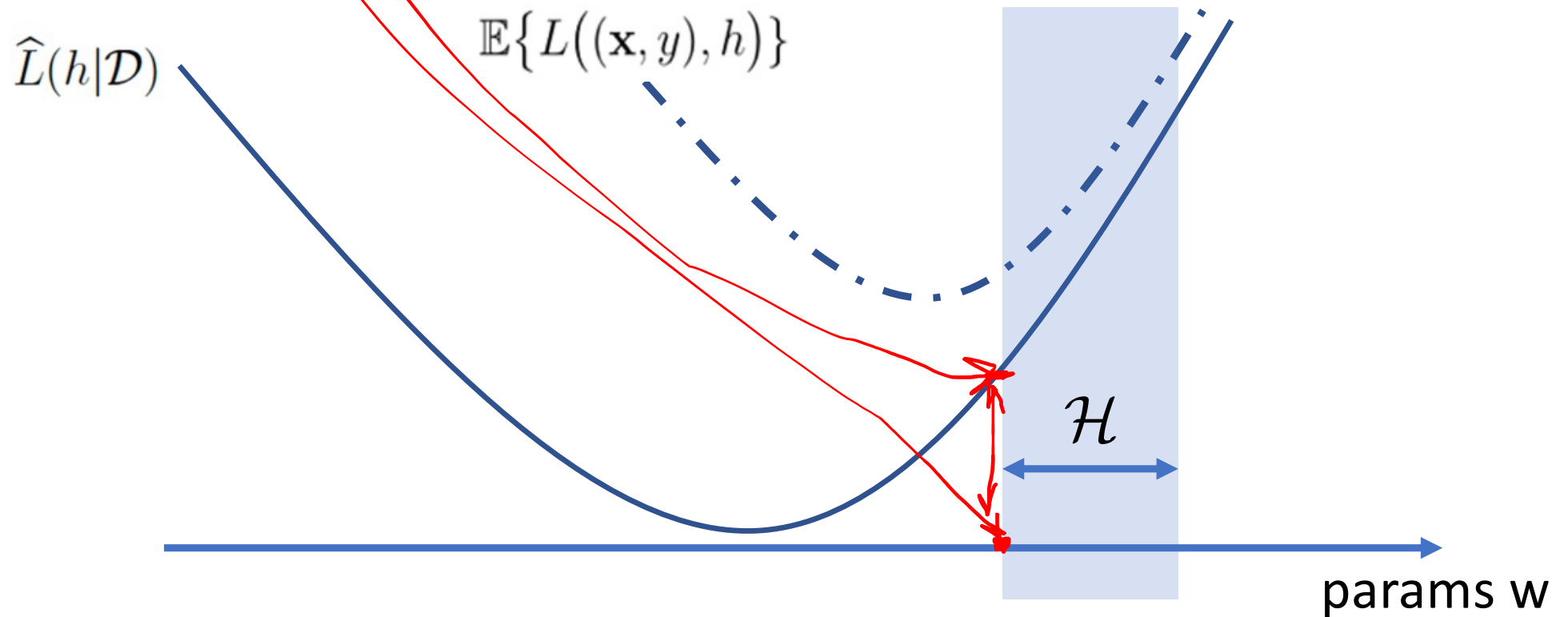
Linear Regression in Python

$$\hat{\mathbf{w}} = \underset{\mathbf{w} \in \mathbb{R}^n}{\operatorname{argmin}} \left(\frac{1}{m} \sum_{i=1}^m \left(y^{(i)} - \mathbf{w}^T \mathbf{x}^{(i)} \right)^2 \right). \quad (4.5)$$

```
In [81]: # Create a linear regression model
         lr = LinearRegression()
         # Fit the model to our data in order to get the coefficients
         lr = lr.fit(features, labels)
```

$$\mathbf{X} = \begin{pmatrix} \mathbf{x}^{(1)} & \dots & \mathbf{x}^{(m)} \end{pmatrix}^T \in \mathbb{R}^{m \times n} \quad \mathbf{y} = \left(y^{(1)}, \dots, y^{(m)} \right)^T \in \mathbb{R}^m$$

```
# create and train a linear model
lr = LinearRegression()
lr = lr.fit(X, y)
w_hat = lr.coef_
trainerr = mean_squared_error(lr.predict(X), y)
```



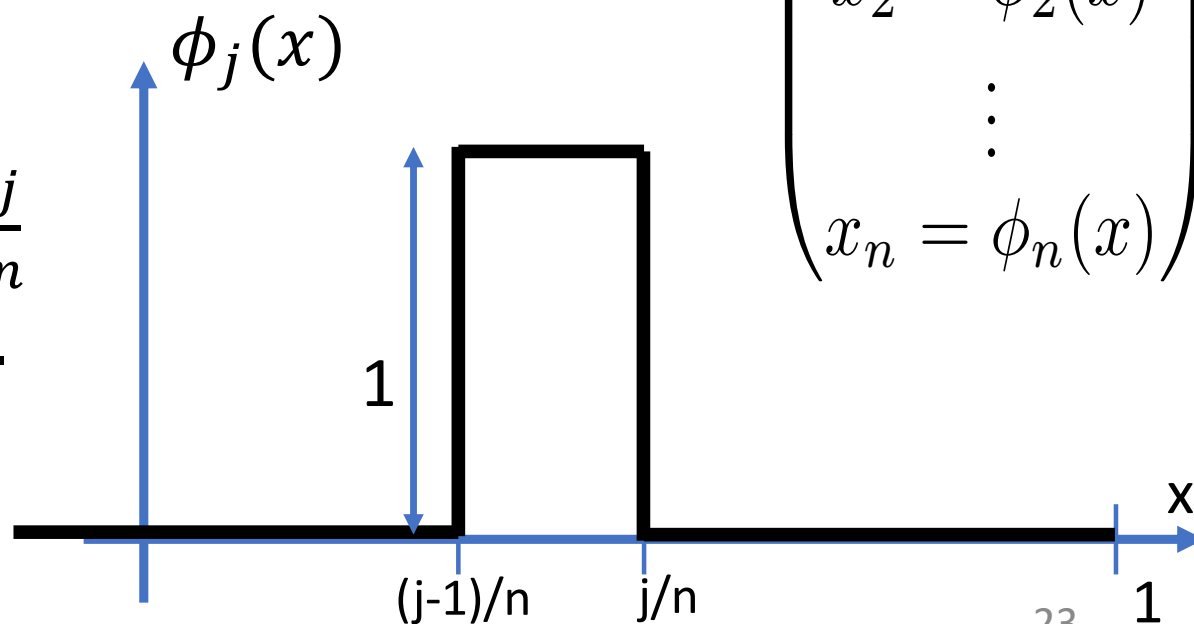
Upgrade Linear Model with new Features !

- consider data points with **single numeric feature x**

- **construct** new features x_1, \dots, x_n

- $x_j = \begin{cases} 1 & \text{for } \frac{j-1}{n} \leq x \leq \frac{j}{n} \\ 0 & \text{for all other } x. \end{cases}$

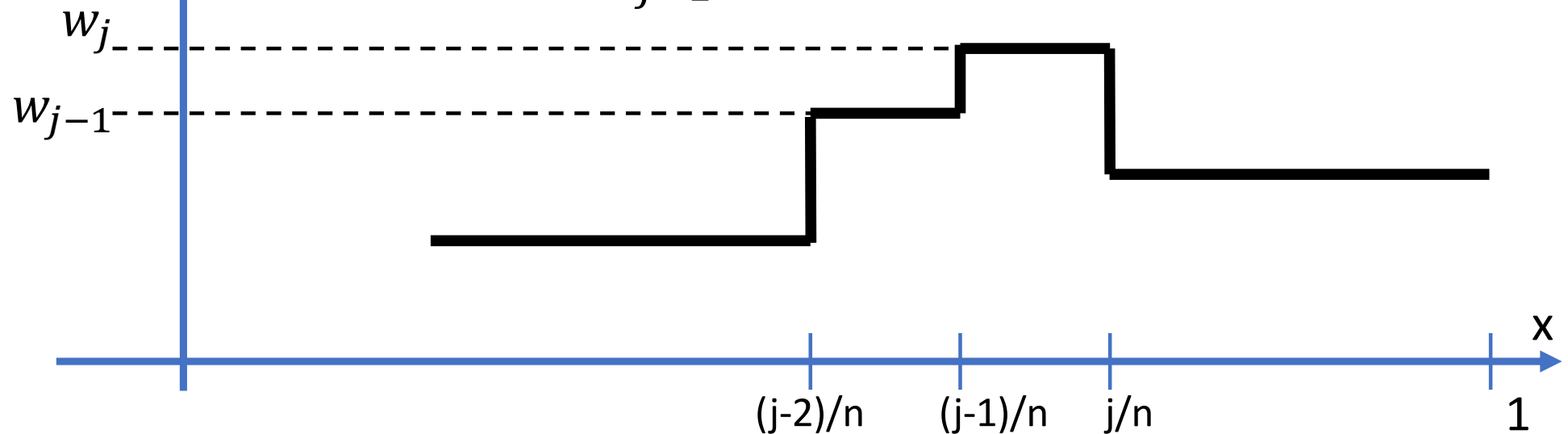
$$\begin{pmatrix} x_1 = \phi_1(x) \\ x_2 = \phi_2(x) \\ \vdots \\ x_n = \phi_n(x) \end{pmatrix}$$



You Can Do Anything with Linear Predictors!

- $h(x)$ is linear in new features but non-linear in raw feature x !

$$h(x_1, \dots, x_n) = \sum_{j=1}^n w_j x_j$$



Polynomial Regression

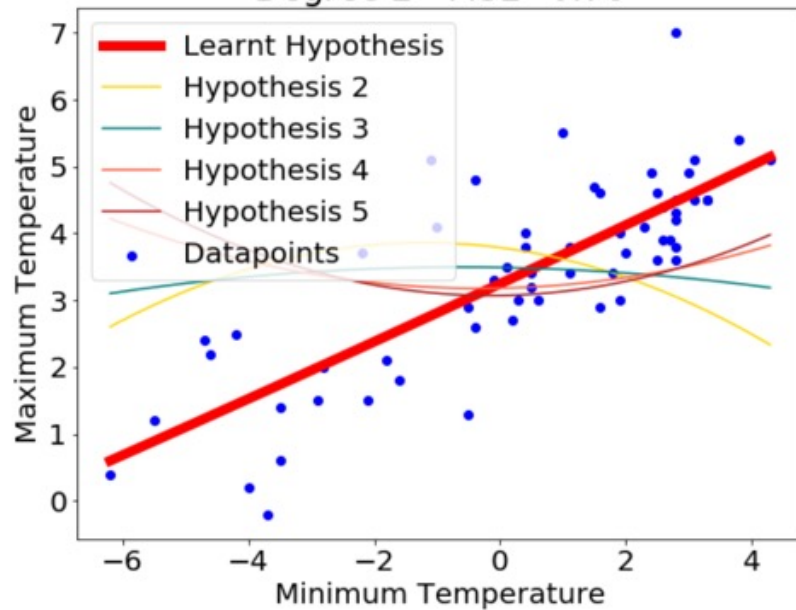
$$\mathcal{H}_{\text{poly}}^{(n)} = \{h^{(\mathbf{w})} : \mathbb{R} \rightarrow \mathbb{R} : h^{(\mathbf{w})}(x) = \sum_{j=1}^n w_j x^{j-1},$$

with some $\mathbf{w} = (w_1, \dots, w_n)^T \in \mathbb{R}^n\}$. (3.4)

Polynomial Regression

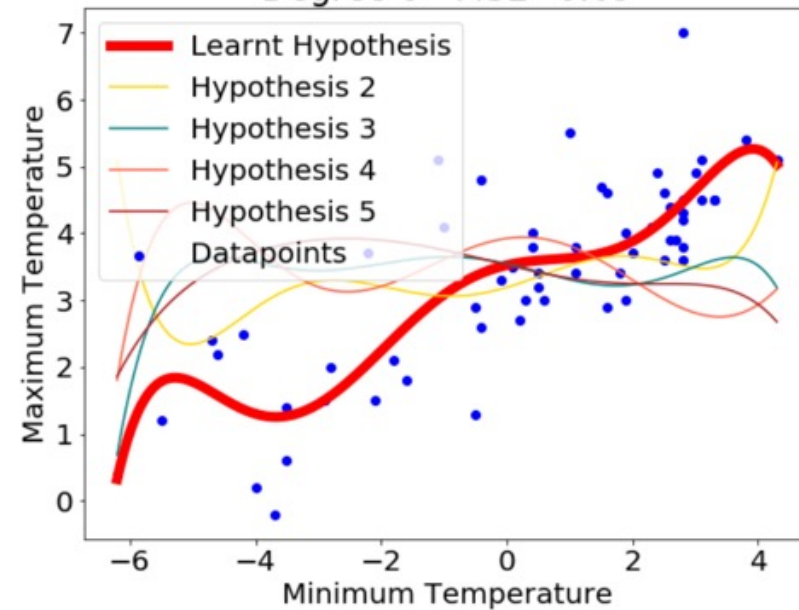
$\mathcal{H}^{(2)}$

Degree 2 - MSE=0.76



$\mathcal{H}^{(6)}$

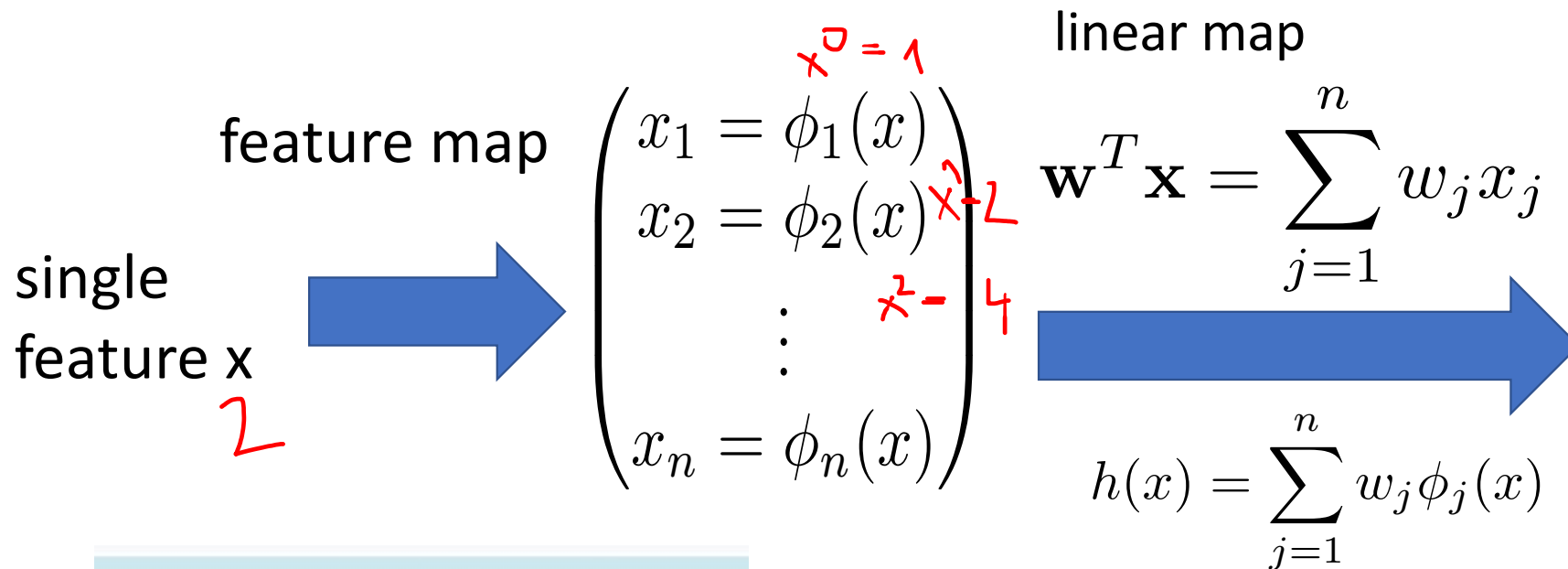
Degree 6 - MSE=0.69



from notebook of TA George

https://github.com/alexjungaalto/cs-c3240spring2022/blob/main/George_Demo_PolynomialRegression.ipynb

Polynomial Regression = Lin. Reg. with Feature Transform.



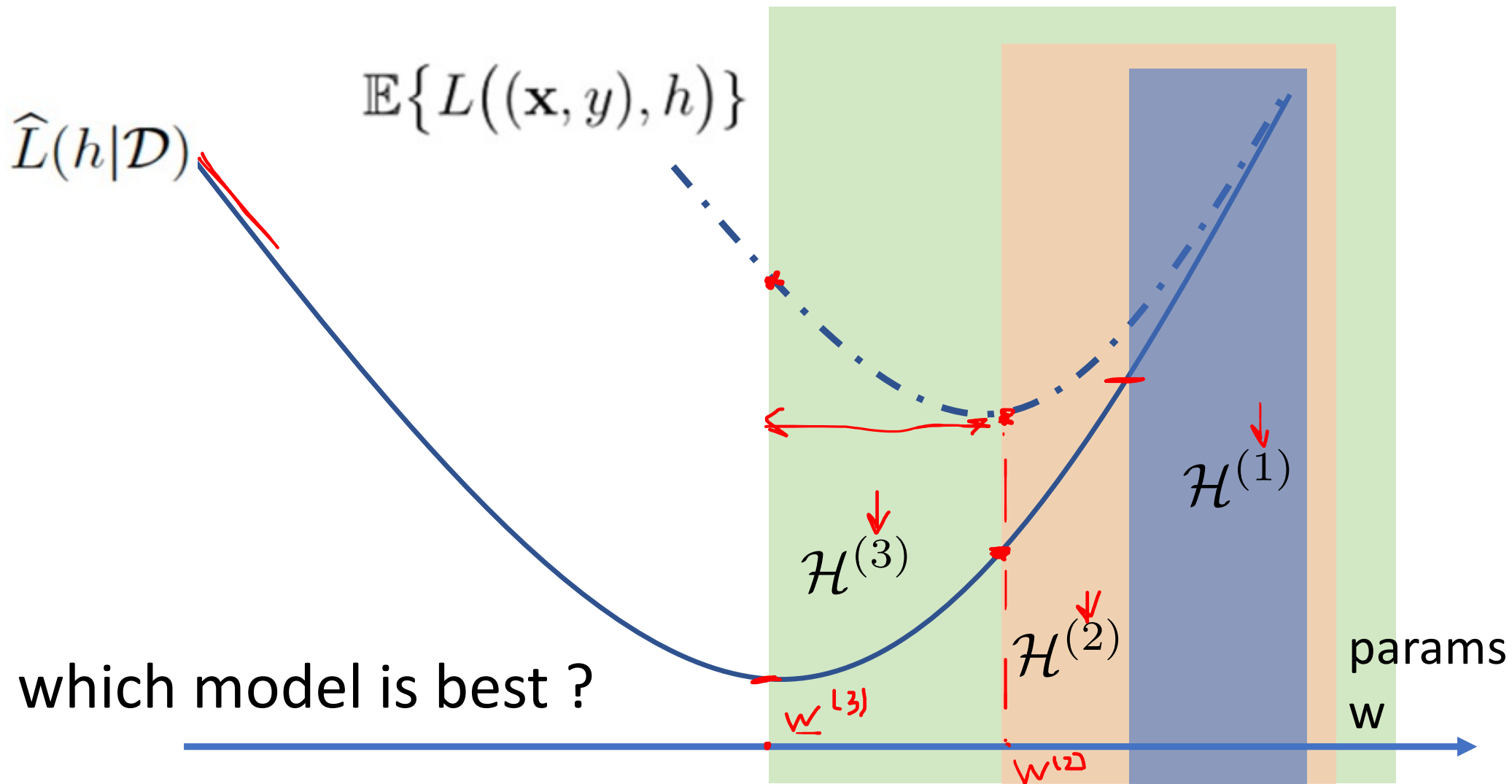
```
sklearn.preprocessing.PolynomialFeatures
.preprocessing.PolynomialFeatures(degree=2, *, interaction_only=False, include_bias=True)
```

```
sklearn.linear_model.LinearRegression
class sklearn.linear_model.LinearRegression(*, fit_intercept=True, normalize='deprecated', copy_X=True, n_positive=False)
```

Polynomial Features

we can use anything as features that can be computed or measured easily !

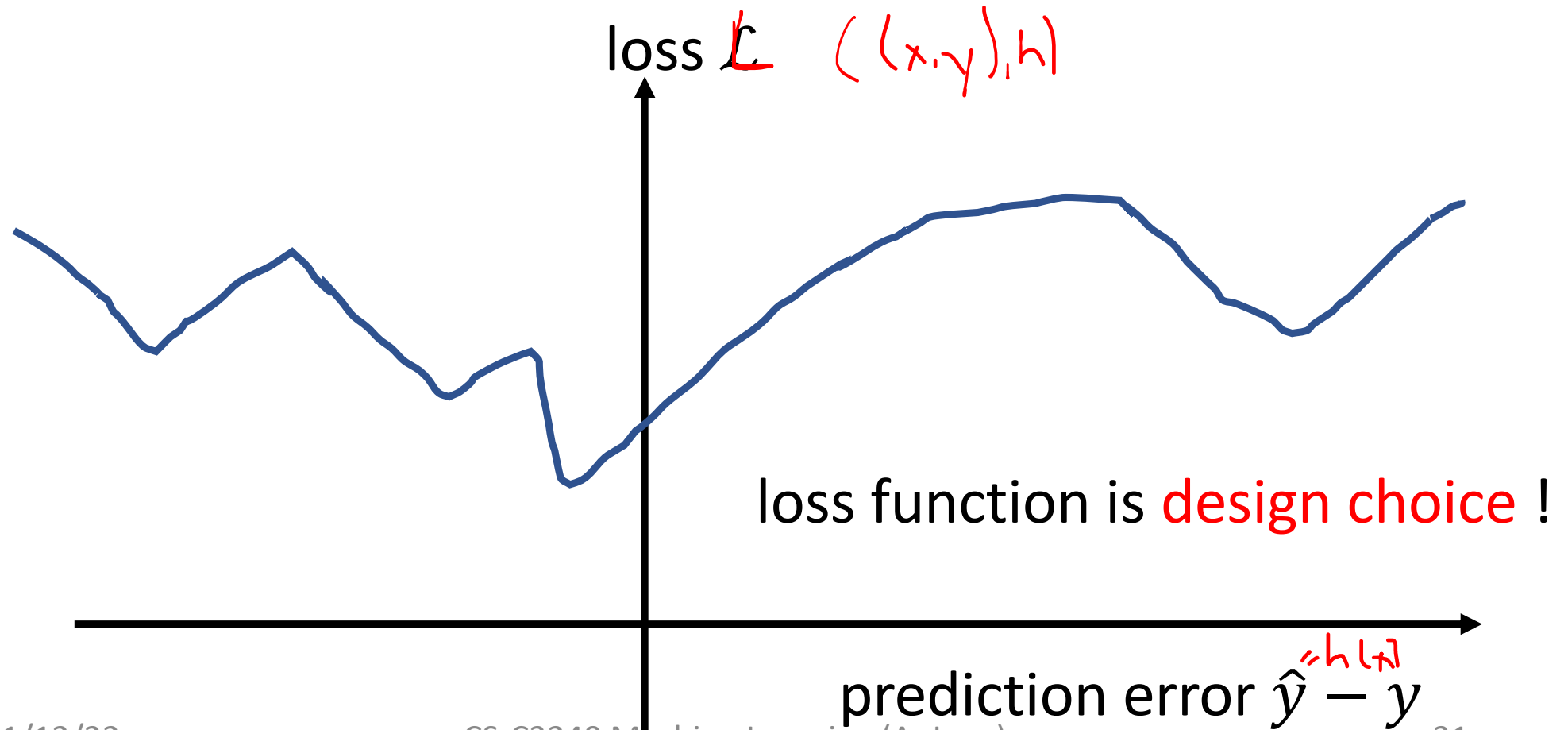
	Date	Max temp	Min temp	(Min temp)^2
0	2020-2-1	3.0	1.9	3.61
1	2020-2-2	4.9 3.0	2.4	5.76
2	2020-2-3	2.6	-0.4	0.16
3	2020-2-4	-0.2	-3.7	13.69
4	2020-2-5	2.5	-4.2	17.64



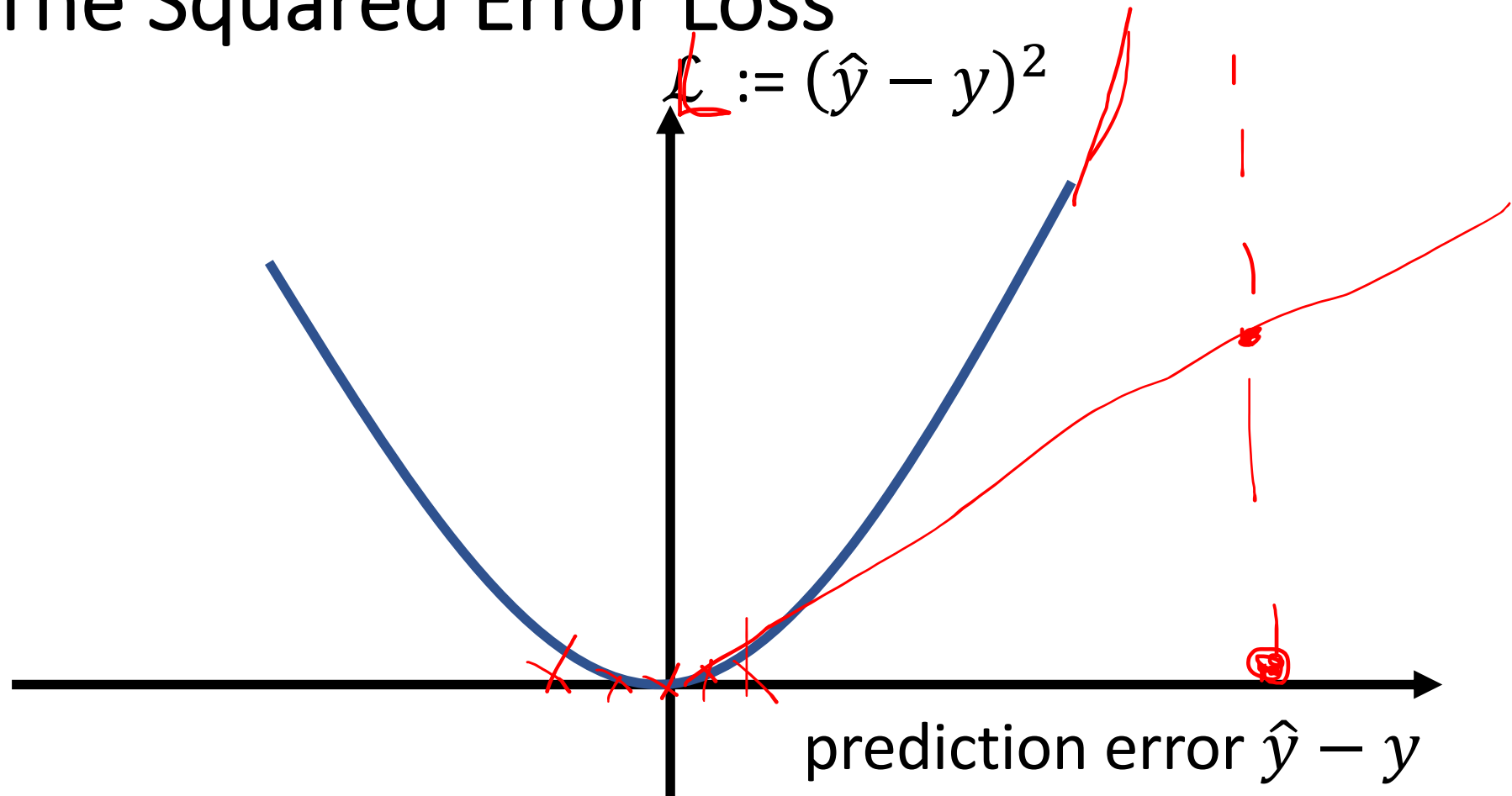
which model is best ?

Design Choice: Loss Function

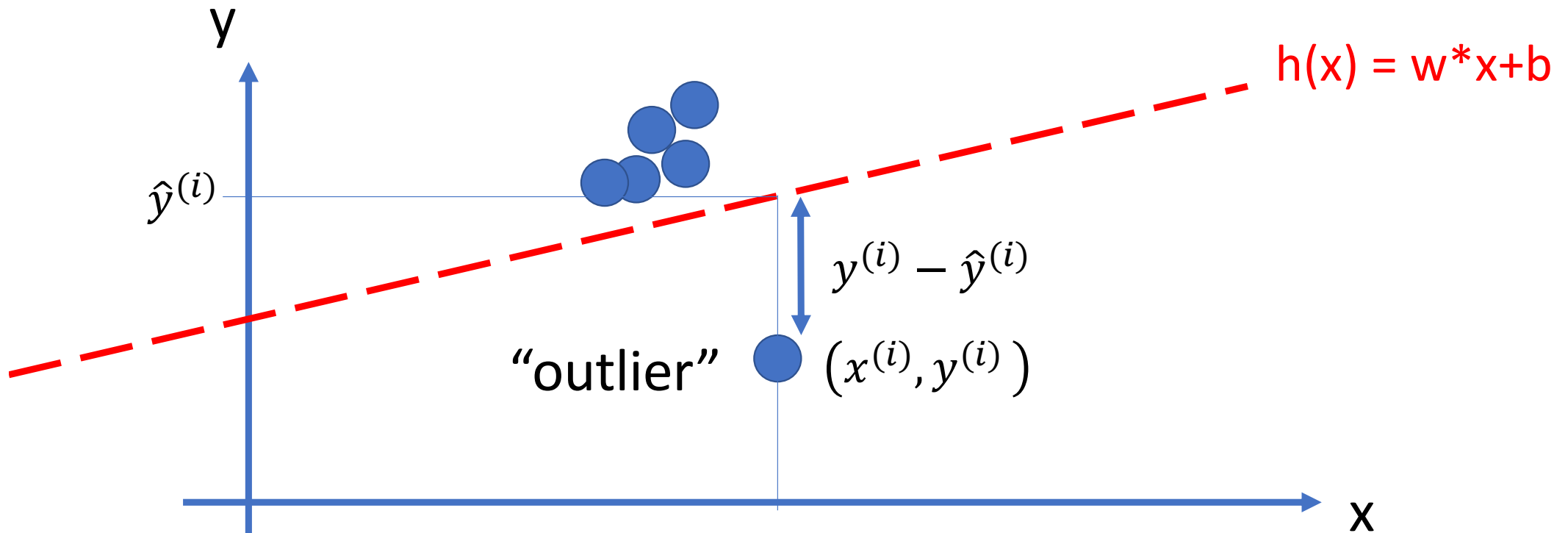
Measuring Error Size via Loss Functions



The Squared Error Loss

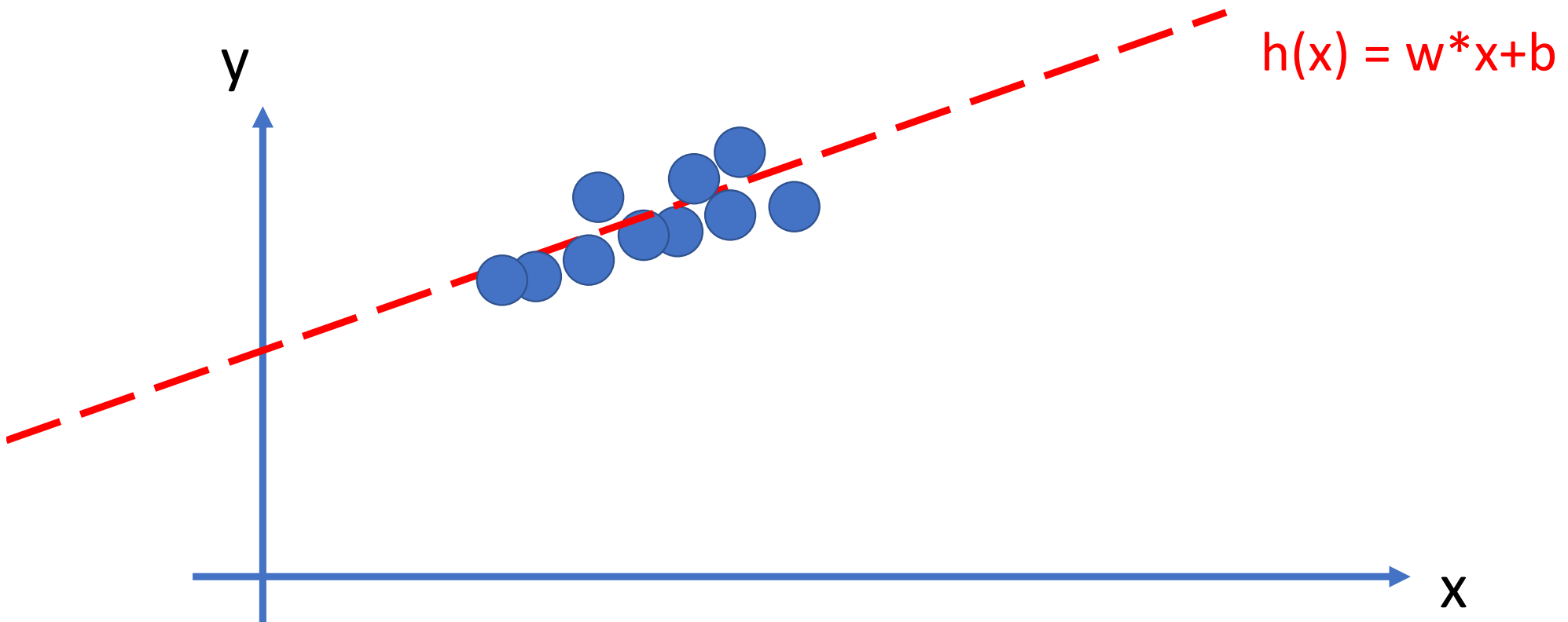


Squared Error Loss Sensitive to Outliers

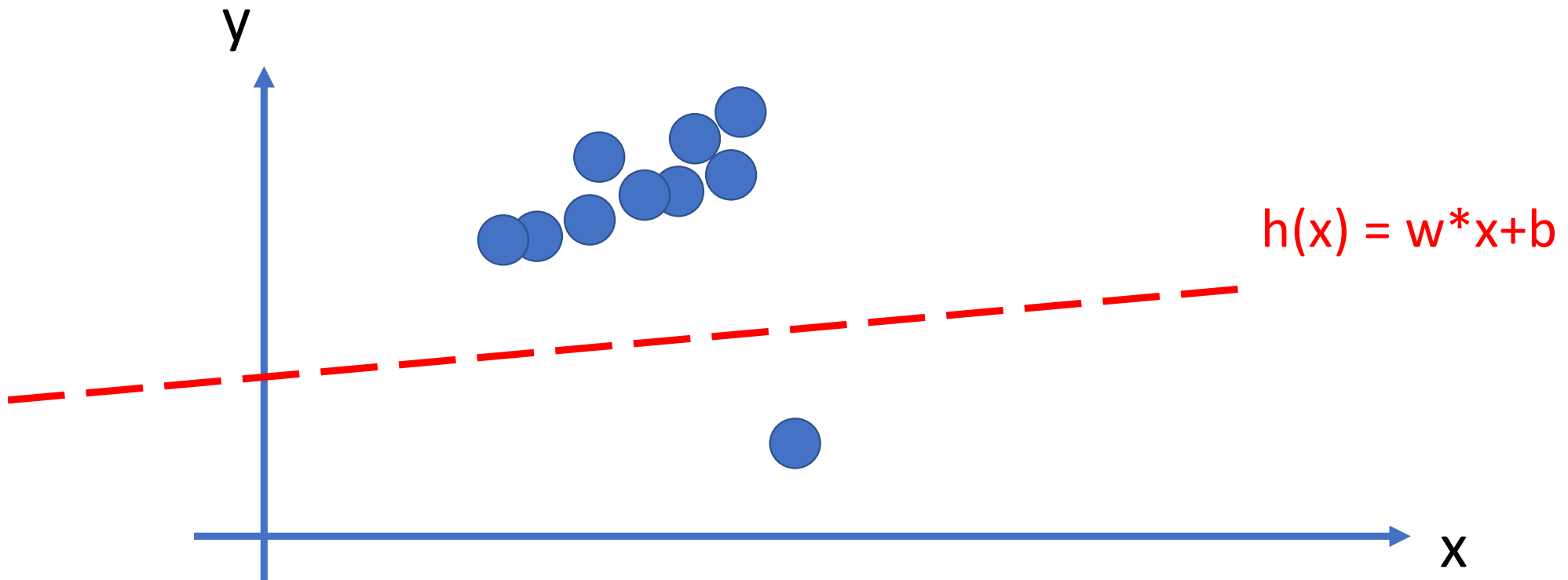


min. squared error loss forces predictor towards outlier

Train Linear Model on “Clean Data”

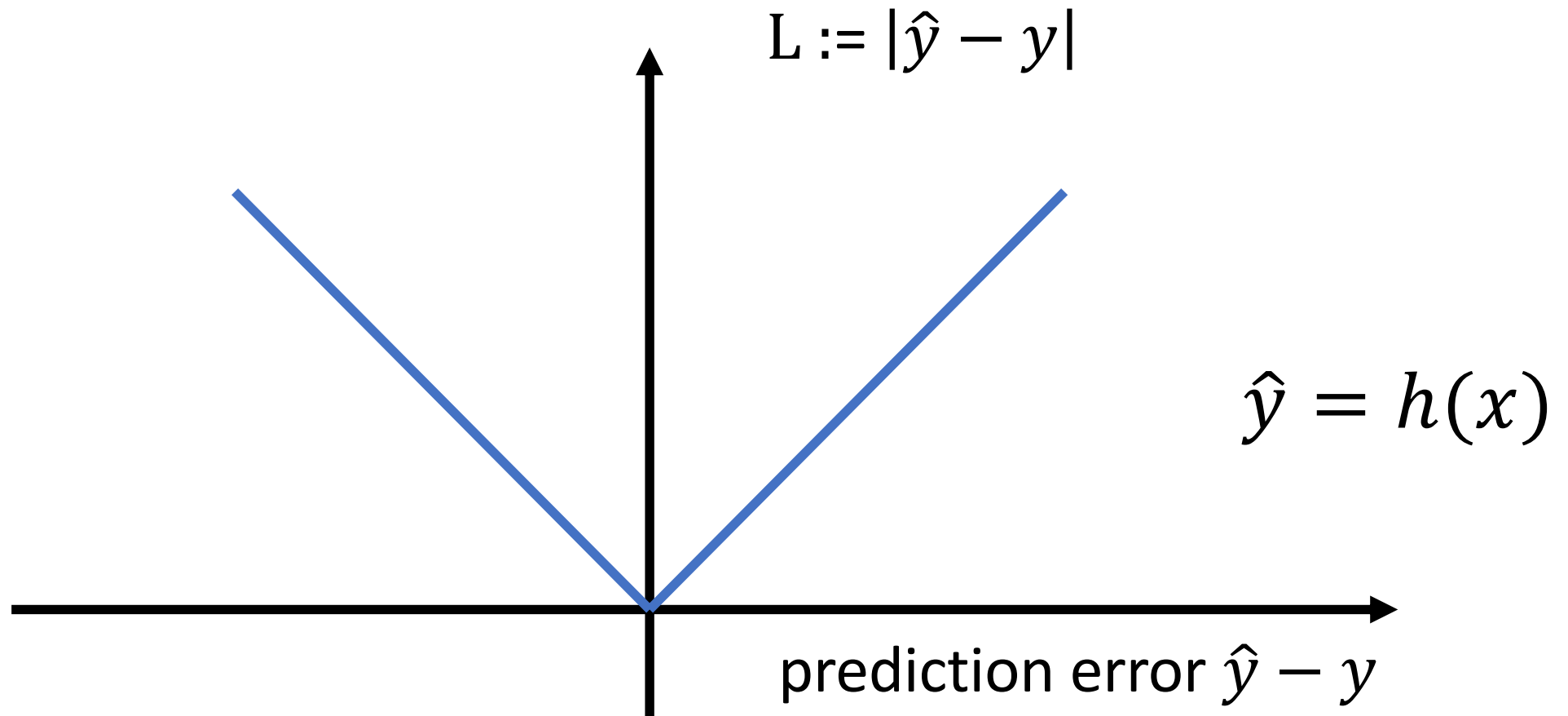


Training Set with a SINGLE OUTLIER !

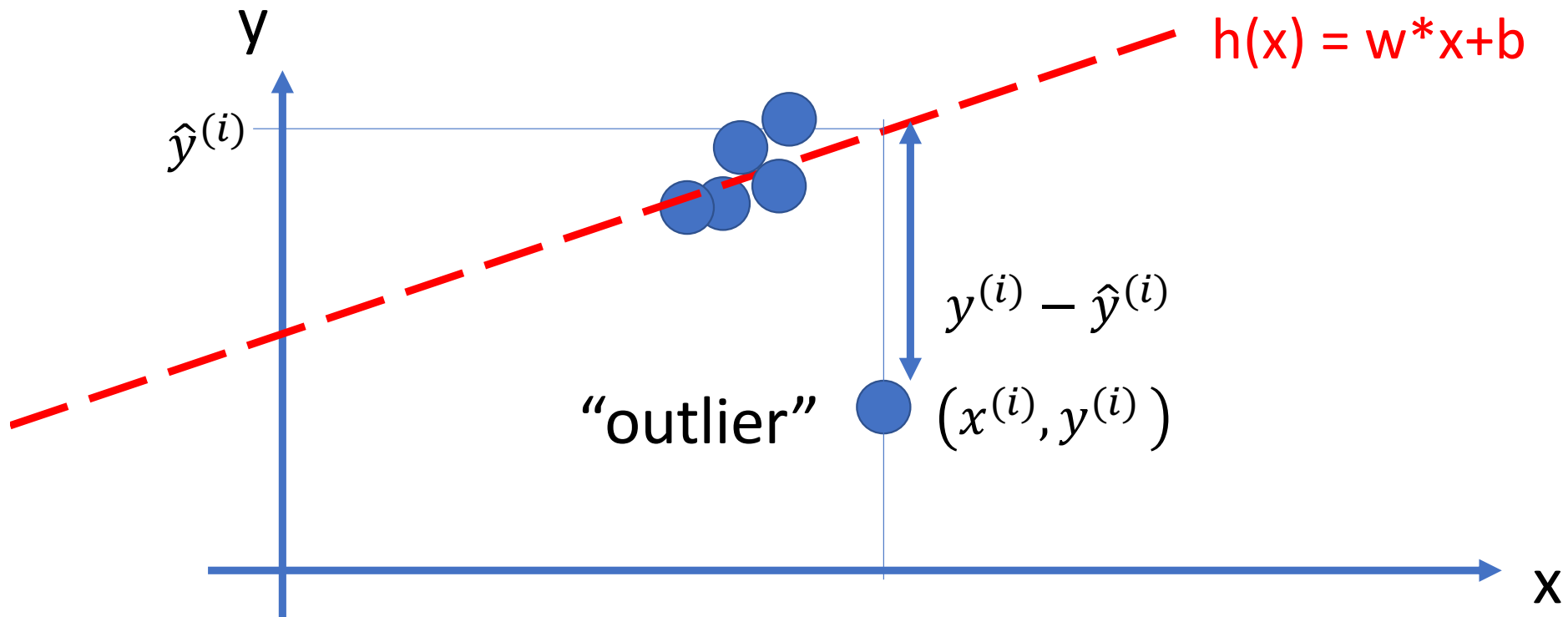


How to make learning robust against presence of few outliers in training set ?

The Absolute Error Loss

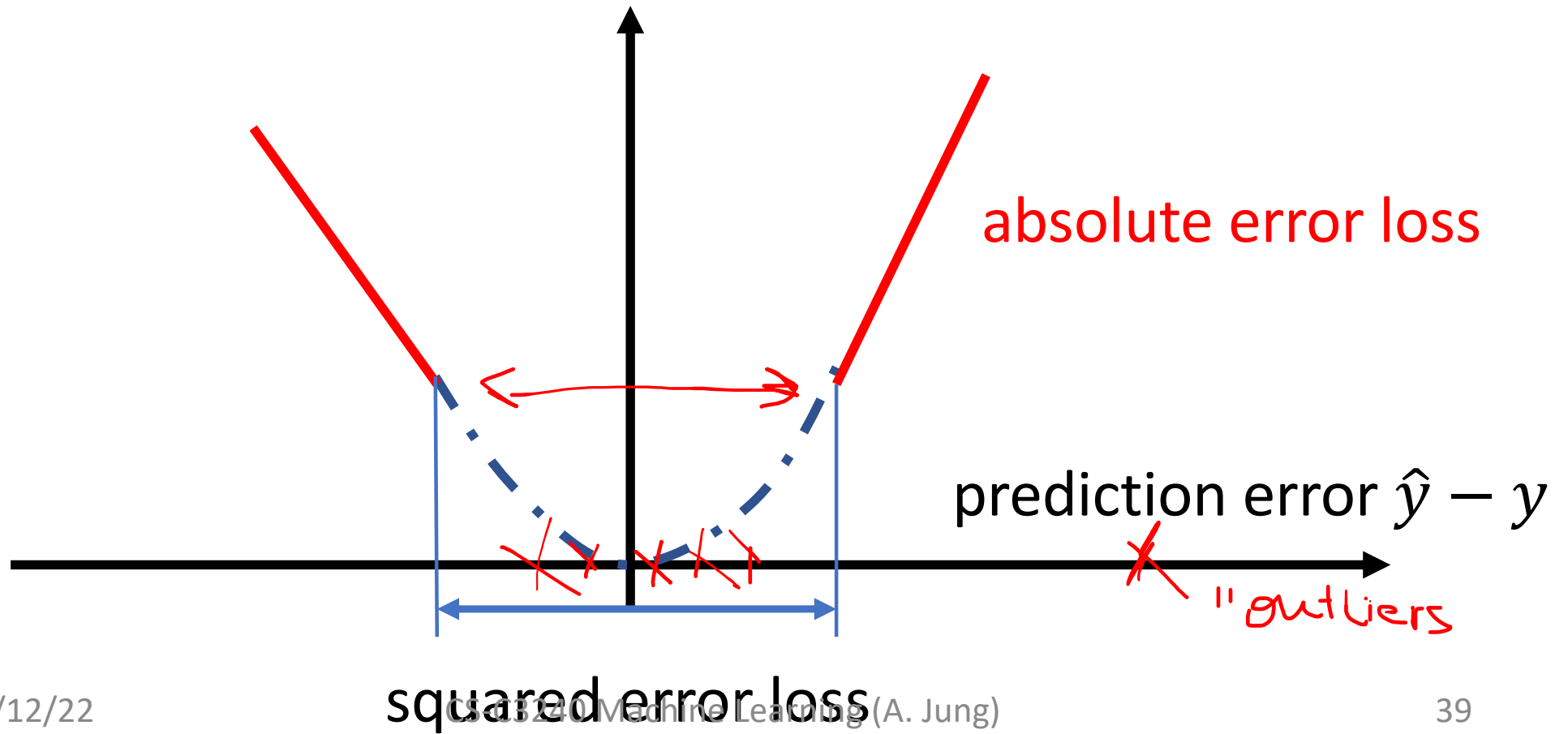


Absolute Error Loss Robust to Outliers

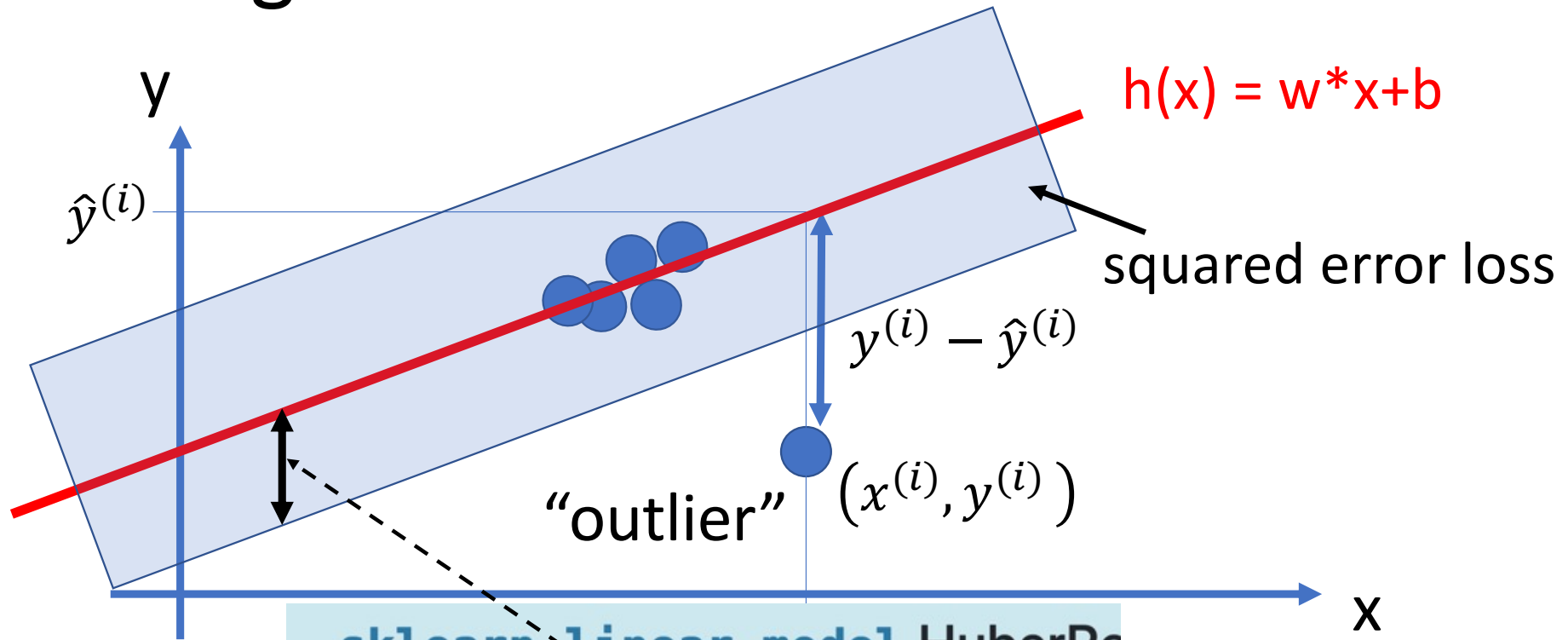


absolute error **"tolerates"** few outliers

Huber Loss

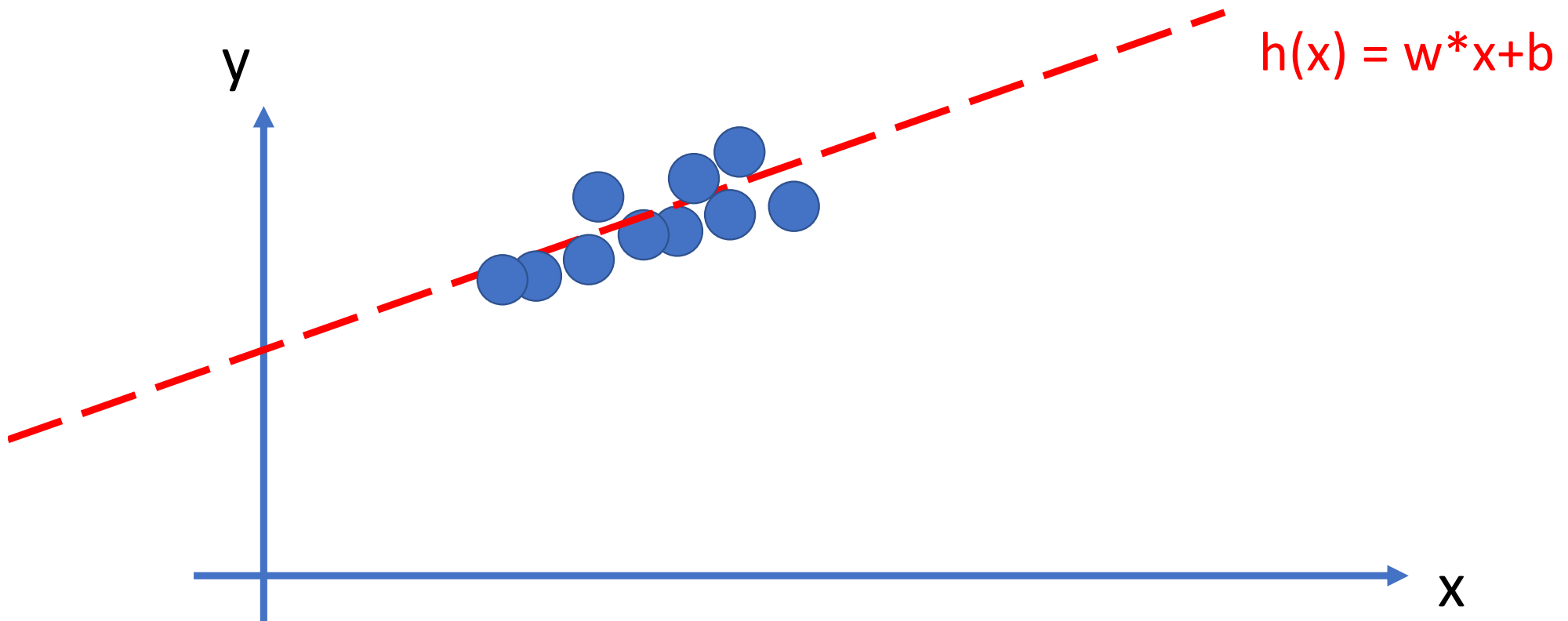


Fitting Linear Predictor with Huber Loss

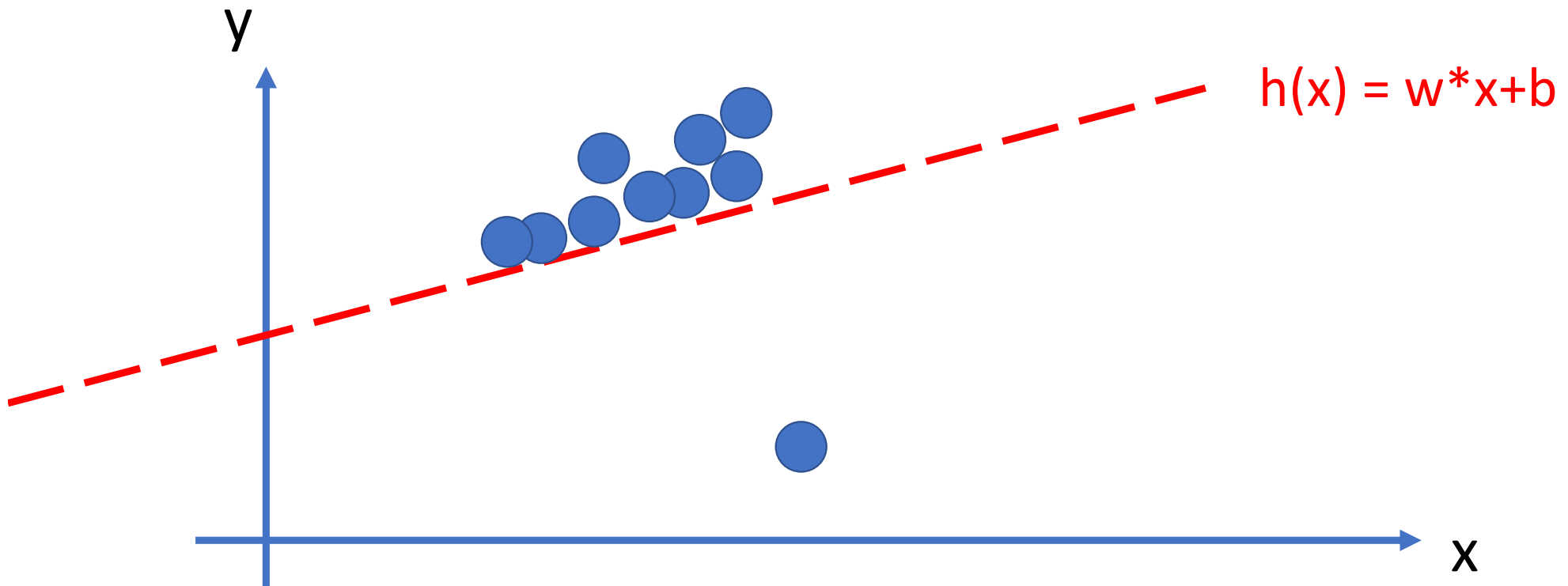


```
sklearn.linear_model.HuberRe  
near_model.HuberRegressor(epsilon=1.35, max_iter=100, alpha=0.0  
tol=1e-05)
```


Train Linear Model on “Clean Data”



Training Set with a SINGLE OUTLIER !



Huber vs. Squared Error Loss

Squared Error

- cvx and diff.able
- minimized via simple gradient descent
- sensitive to outliers

Huber

- cvx and non-diff.
- requires more advanced opt. methods
- robust against outliers

Summary

- ultimate quality measure: expected loss or risk
- approximate risk by average loss (empirical risk)
- many ML methods are instances of ERM
- three design choices of ERM: data, model and loss
- ERM can fail if empirical risk deviates from risk

What's Next ?

- work on assignment A1 on <https://jupyter.cs.aalto.fi/>
- review Stage 1 of ML Project (Section “ML Project”)
- next Lecture on Mo. 17.01.2022 on Classification