

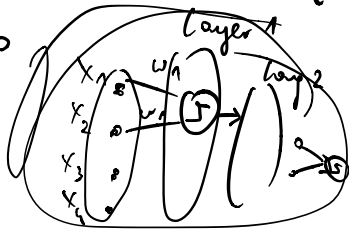
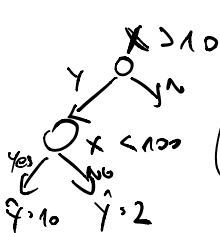
Machine Learning

Data

- Images
feature $\underline{x} \in \mathbb{R}^n$

label $y = \# \text{ persons on image}$

$y = \begin{cases} > 5 \text{ persons} \\ \emptyset \text{ persons} \end{cases}$



Hypothesis Space



$$h(\cdot) : \underline{x} \mapsto \hat{y} \in \mathbb{R}$$

"hypothesis"
"predictor"

"classifier"

$\{0, 1\}$
"val"
"deg"

$$h(\underline{x}) \leq \underline{w}^T \underline{x}$$

Loss

Squared error

$$(\hat{y} - y)^2$$

absolute error

$$-|\hat{y} - y|$$

$$\mathbb{I}_{\hat{y} \neq y} = \begin{cases} 1 & \hat{y} \neq y \\ 0 & \hat{y} = y \end{cases}$$

"0/1 loss"
 $\hat{y} \neq y$
 $\emptyset \quad y = \hat{y}$

$$(\underline{x}^{(1)}, y^{(1)})$$

$$(\underline{x}^{(2)}, y^{(2)})$$

$$\dots (\underline{x}^{(m)}, y^{(m)})$$

$$\underline{x}^{(m)} = \begin{pmatrix} x_1^{(m)} \\ x_2^{(m)} \\ \vdots \\ x_n^{(m)} \end{pmatrix}$$

linear hypothesis space

$$h^{(w)}(\underline{x}) = \underline{w}^T \underline{x} = \sum_{r=1}^n w_r x_r$$

Squared error loss

$\underline{w} \dots$ "weight vector"

"Linear Regression"

best linear predictor?

$$\text{minimize } \frac{1}{m} \sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)})^2$$

true label

predicted label

$$\min_{\underline{w} \in \mathbb{R}^n} E(\underline{w}) = \frac{1}{m} \sum_{i=1}^m (y^{(i)} - \underline{w}^T \underline{x}^{(i)})^2$$

$$\approx E^{(val)} - E^{(train)}$$

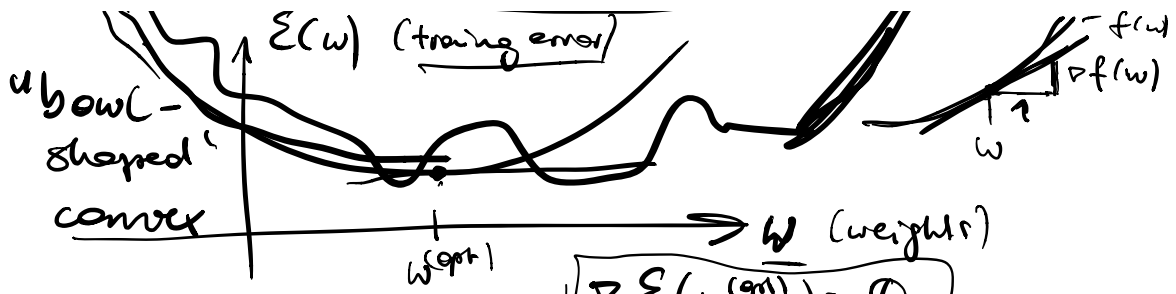
complexity

$$R(\underline{w})$$

training error

prediction

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



$$\nabla E(\underline{w}^{(opt)}) = \underline{0}$$

label vector $\underline{y} = \begin{pmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{pmatrix} \in \mathbb{R}^m$

feature matrix $\underline{X} = \begin{pmatrix} (x^{(1)})^T \\ \vdots \\ (x^{(m)})^T \end{pmatrix} \in \mathbb{R}^{m \times n}$ $\underline{X} = \begin{pmatrix} x_1 \\ \vdots \\ x_m \end{pmatrix}$

$$E(w) = \frac{1}{m} \|\underline{X} - \underline{X}\underline{w}\|_2^2$$

$$\nabla E(\underline{w}^{(opt)}) = \underline{0}$$

$$\|\underline{X}\|_2^2 = \sum_{r=1}^n x_r^2$$

$$n > m \Rightarrow E(w) > 0$$

$$(-2/m) \underline{X}^T (\underline{y} - \underline{X}\underline{w}^{(opt)}) \neq \underline{0} \quad \text{OVER FITTING!}$$

$$\cdot \underline{C}^{-1} \quad \left[(2/m) \underline{X}^T \underline{X} \underline{w}^{(opt)} = (+2/m) \underline{X}^T \underline{y} \right]$$

$$\underline{w}^{(opt)} = \underline{C}^{-1} \underline{X}^T \underline{y}$$

$\underline{C} = \underline{X}^T \underline{X}$ non-invertible / singular
 $n > m$

training $\rightarrow \underline{w}^{(opt)} \rightarrow$ validate!

$$E^{(val)}(\underline{w}^{(opt)}) = \sum_{i=m+1}^{m+m_{val}} (\underline{w}^T \underline{x}^{(i)} - y^{(i)})^2$$

$$E \approx E^{(val)} \approx E^{reference}$$

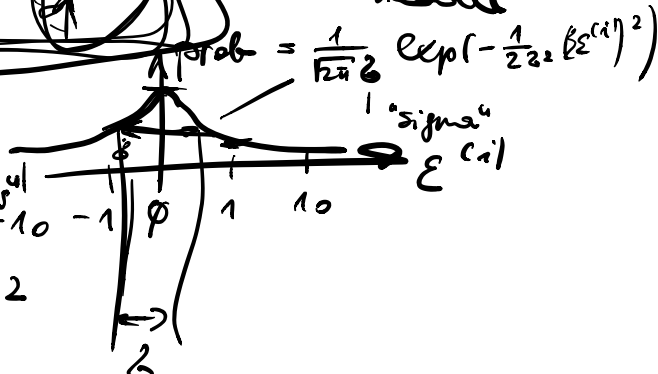
Observation model

$$Y^{(i)} = \underline{w}_0^T \underline{X}^{(i)} + \varepsilon^{(i)}$$

$$h(\underline{x}) = \hat{y}$$

$$\underline{w}^T \underline{x}^{(i)} = \hat{y}^{(i)} \neq y^{(i)}$$

$$\frac{1}{m} \sum (y^{(i)} - \hat{y}^{(i)})^2 \leq \sigma^2$$



$$\begin{pmatrix} \underline{x}^{(1)} & y^{(1)} \end{pmatrix} \dots \begin{pmatrix} \underline{x}^{(m)} & y^{(m)} \end{pmatrix}$$

$$\begin{pmatrix} y^{(1)} \end{pmatrix} = \begin{pmatrix} \underline{w}_0^T \underline{x}^{(1)} \end{pmatrix} + \begin{pmatrix} \varepsilon^{(1)} \end{pmatrix} \dots y^{(m)} = \begin{pmatrix} \underline{w}_0^T \underline{x}^{(m)} \end{pmatrix} + \begin{pmatrix} \varepsilon^{(m)} \end{pmatrix}$$

$\varepsilon^{(i)} \in \mathbb{R}$

$$\underline{x}^{(1)} = \underline{x}^{(2)} = \dots = \underline{x}^{(m)}$$

empirical mean

$$\hat{\mu} = \frac{1}{m} \sum_{i=1}^m y^{(i)}$$

$$\hat{\sigma}^2 = \frac{1}{m} \sum_{i=1}^m (y^{(i)} - \hat{\mu})^2 \quad (\approx \sigma^2)$$

Optimizer

SGD

Descent

Stochastic Gradient

$$\min_{\underline{w}} \frac{1}{m} \sum (y^{(i)} - \underline{w}^T \underline{x}^{(i)})^2 = \min_{\underline{w}} \frac{1}{m} \|\underline{y} - \underline{X} \underline{w}\|_2^2$$

optimal weight vector

$$\nabla \mathcal{E}(\underline{w}^{(opt)}) = \emptyset$$

$\mathcal{E}(\underline{w})$ training error



Gradient Descent

initial guess $\underline{w}^{(0)} \rightarrow \underline{w}^{(1)} \rightarrow \underline{w}^{(2)}$

α too large
 α too small

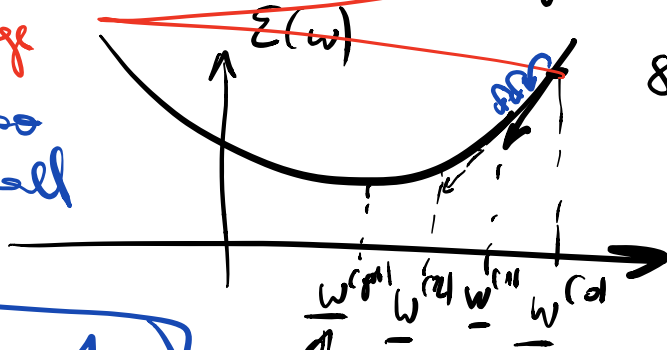
$\underline{w}^{(k+1)}$
new

$\underline{w}^{(k)}$
old guess

$\nabla \mathcal{E}(\underline{w}^{(k)})$

gradient at old guess

step-size



"partial" derivative
 $\frac{\partial \mathcal{E}(\underline{w})}{\partial \underline{w}}$

$\alpha = 1$

optimal

$$\nabla \mathcal{E}(\underline{w}) = 2 \underline{w} \underline{X}^T (\underline{y} - \underline{X} \underline{w})$$

linear regression

Squared error loss +
linear predictor

$$\nabla \mathcal{E}(\underline{w}^{(k)}) = \text{error at } (2 \underline{w}^{(k)}) \underline{X}^T (\underline{y} - \underline{X} \underline{w}^{(k)}) \text{ iterate}$$

$$\begin{pmatrix} y^{(1)} - \hat{y}^{(1)} \\ \vdots \\ y^{(n)} - \hat{y}^{(n)} \end{pmatrix}$$