# Machine Learning

## Linear Models

Fabio Vandin                    October 28$^{th}$, 2022

# Linear Regression $\neq$ Regression

$\mathcal{X} = \mathbb{R}^d$, $\mathcal{Y} = \mathbb{R}$

# Linear Regression

$\mathcal{X} = \mathbb{R}^d$, $\mathcal{Y} = \mathbb{R}$

Hypothesis class:

$$\mathcal{H}_{reg} = L_d = \{\mathbf{x} \to \langle \mathbf{w}, \mathbf{x} \rangle + b : \mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}\}$$

# Linear Regression

$\mathcal{X} = \mathbb{R}^d$, $\mathcal{Y} = \mathbb{R}$

Hypothesis class:

$$\mathcal{H}_{reg} = L_d = \{\mathbf{x} \to \langle \mathbf{w}, \mathbf{x} \rangle + b : \mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}\}$$

**Note:** $h \in \mathcal{H}_{reg} : \mathbb{R}^d \to \mathbb{R}$

Commonly used loss function: *squared-loss*

$$\ell(h, (\mathbf{x}, y)) \stackrel{def}{=} (h(\mathbf{x}) - y)^2$$

ERM for regression with linear models and squared loss

# Linear Regression

$\mathcal{X} = \mathbb{R}^d$, $\mathcal{Y} = \mathbb{R}$

Hypothesis class:

$$\mathcal{H}_{reg} = L_d = \{\mathbf{x} \to \langle \mathbf{w}, \mathbf{x} \rangle + b : \mathbf{w} \in \mathbb{R}^d, b \in \mathbb{R}\}$$

**Note:** $h \in \mathcal{H}_{reg} : \mathbb{R}^d \to \mathbb{R}$

Commonly used loss function: *squared-loss*

$$\ell(h, (\mathbf{x}, y)) \overset{def}{=} (h(\mathbf{x}) - y)^2$$

$\Rightarrow$ empirical risk function (training error): *Mean Squared Error*

$$S = \left\{ (\vec{x}_1, y_1), \dots, (\vec{x}_m, y_m) \right\} \quad L_S(h) = \frac{1}{m} \sum_{i=1}^{m} \underbrace{(h(\mathbf{x}_i) - y_i)^2}_{\ell\left(h, (\vec{x}_i, y_i)\right)}$$
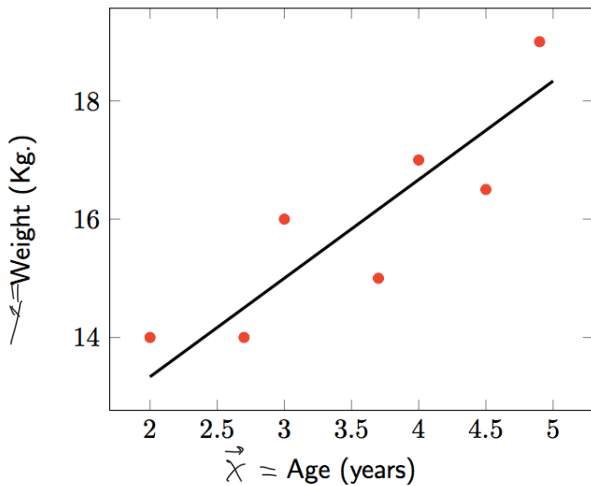
2

# Linear Regression - Example



$d = 1$

training set

line

$h_{\vec{w}}(\vec{x}) = w_1 x_1 + w_0$

$\vec{x} = [x_1] \rightarrow \begin{bmatrix} 1 \\ x_1 \end{bmatrix}$

$[x_1] = \vec{x} \in \mathbb{R}$

3

# Linear Regression - Example

$d = 1$

# Least Squares

How to find a ERM hypothesis? *Least Squares* algorithm

Best hypothesis:

here: $\vec{x_i} = \begin{bmatrix} 1, & x_{i1}, & \cdots, & x_{id} \end{bmatrix}^{\top}$

$\vec{w} = \begin{bmatrix} w_o, & w_1, & \cdots, & w_d \end{bmatrix}^{\top}$

$$\arg\min_{\mathbf{w}} L_S(h_{\mathbf{w}}) = \arg\min_{\mathbf{w}} \frac{1}{m} \sum_{i=1}^{m} (\langle \mathbf{w}, \mathbf{x}_i \rangle - y_i)^2$$

# Least Squares

How to find a ERM hypothesis? *Least Squares* algorithm

Best hypothesis:

$$\arg\min_{\mathbf{w}} L_S(h_{\mathbf{w}}) = \arg\min_{\mathbf{w}} \frac{1}{m} \sum_{i=1}^{m} (\langle \mathbf{w}, \mathbf{x}_i \rangle - y_i)^2$$

Equivalent formulation: **w** minimizing *Residual Sum of Squares* (RSS), i.e.

$$\arg\min_{\mathbf{w}} \underbrace{\sum_{i=1}^{m} (\langle \mathbf{w}, \mathbf{x}_i \rangle - y_i)^2}_{RSS}$$

# RSS: Matrix Form

Let



$$\mathbf{X} = \begin{bmatrix} \cdots & \mathbf{x}_1 & \cdots \\ \cdots & \mathbf{x}_2 & \cdots \\ \cdots & \vdots & \cdots \\ \cdots & \mathbf{x}_m & \cdots \end{bmatrix}$$

$\mathbf{X}$: *design matrix*

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}$$

$$S = \left\{ \left( \vec{x}_1, y_1 \right), \left( \vec{x}_2, y_2 \right), \cdots, \left( \vec{x}_m, y_m \right) \right\}$$

# RSS: Matrix Form

Let

$$\mathbf{X} = \begin{bmatrix} \cdots & \mathbf{x}_1 & \cdots \\ \cdots & \mathbf{x}_2 & \cdots \\ \cdots & \vdots & \cdots \\ \cdots & \mathbf{x}_m & \cdots \end{bmatrix}$$

$\mathbf{X}$: *design matrix*

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}$$
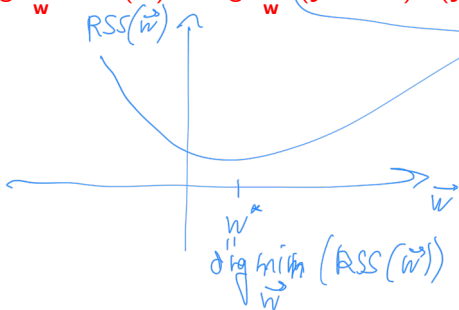
$\Rightarrow$ we have that RSS is

$$\sum_{i=1}^{m} (\langle \mathbf{w}, \mathbf{x}_i \rangle - y_i)^2 = (\mathbf{y} - \mathbf{Xw})^T (\mathbf{y} - \mathbf{Xw})$$

HW: check that the above is true

5

Want to find **w** that minimizes RSS (=*objective function*):

$$\arg\min_{\mathbf{w}} RSS(\mathbf{w}) = \arg\min_{\mathbf{w}} (\mathbf{y} - \mathbf{Xw})^T (\mathbf{y} - \mathbf{Xw})$$



$\vec{w} \in \mathbb{R}$

$RSS(\vec{w})$

$\vec{w}$

$w^*$

$= \arg\min_{\vec{w}} \left( RSS(\vec{w}) \right)$

convex
function

Want to find **w** that minimizes RSS (=*objective function*):

$$\arg\min_{\mathbf{w}} RSS(\mathbf{w}) = \arg\min_{\mathbf{w}} (\mathbf{y} - \mathbf{Xw})^T (\mathbf{y} - \mathbf{Xw})$$

How?

Compute gradient $\frac{\partial RSS(\mathbf{w})}{\partial \mathbf{w}}$ of objective function w.r.t **w** and compare it to 0.

$$\frac{\partial RSS(\mathbf{w})}{\partial \mathbf{w}} = -2\mathbf{X}^T(\mathbf{y} - \mathbf{Xw})$$

Then we need to find **w** such that

$$-2\mathbf{X}^T(\mathbf{y} - \mathbf{Xw}) = 0$$

$$-2\mathbf{X}^{T}(\mathbf{y} - \mathbf{Xw}) = 0$$

$$-2X^{T}\vec{y} + 2X^{T}X\vec{w} = 0$$

$$2X^{T}X\vec{w} = 2X^{T}\vec{y}$$

$$X^{T}X\vec{w} = X^{T}\vec{y}$$

$$A\vec{w} = b$$

$$\vec{w} = \ldots$$

$$A^{-1}A\vec{w} = A^{-1}b$$

$$\vec{w} = A^{-1}b$$

$$\left(X^{T}X\right)^{-1}\left(X^{T}X\right)\vec{w} = \left(X^{T}X\right)^{-1}X^{T}\vec{y}$$

$$\vec{w} = \left(X^{T}X\right)^{-1}X^{T}\vec{y}$$

7

$$-2\mathbf{X}^T(\mathbf{y} - \mathbf{X}\mathbf{w}) = 0$$

is equivalent to

$$\mathbf{X}^T\mathbf{X}\mathbf{w} = \mathbf{X}^T\mathbf{y}$$

If $\mathbf{X}^T\mathbf{X}$ is invertible $\Rightarrow$ solution to ERM problem is:

$$\mathbf{w} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}$$

# Complexity Considerations

We need to compute

$$(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}$$

$$\vec{X} = \begin{bmatrix} 1 & x_1 & \dots & x_d \end{bmatrix}$$

$$S = \left\{ (\vec{x_1}, y_1), \dots, (\vec{x_m}, y_m) \right\}$$

Algorithm:

1. compute $\mathbf{X}^T\mathbf{X}$: product of $(d+1) \times m$ matrix and $m \times (d+1)$ matrix

# Complexity Considerations

We need to compute
$$(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}$$

Algorithm:

1. compute $\mathbf{X}^T\mathbf{X}$: product of $(d+1) \times m$ matrix and $m \times (d+1)$ matrix
2. compute $(\mathbf{X}^T\mathbf{X})^{-1}$ inversion of $(d+1) \times (d+1)$ matrix
3. compute $(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T$: product of $(d+1) \times (d+1)$ matrix and $(d+1) \times m$ matrix
4. compute $(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}$: product of $(d+1) \times m$ matrix and $m \times 1$ matrix

Most expensive operation? Inversion!

$\Rightarrow$ done for $(d+1) \times (d+1)$ matrix

# $\mathbf{X}^T\mathbf{X}$ not invertible?

How do we get $\mathbf{w}$ such that

$$\mathbf{X}^T\mathbf{X}\mathbf{w} = \mathbf{X}^T\mathbf{y}$$

if $\mathbf{X}^T\mathbf{X}$ is not invertible?
Let

$$\mathbf{A} = \mathbf{X}^T\mathbf{X}$$

Let $\mathbf{A}^+$ be the *generalized inverse* of $\mathbf{A}$, i.e.:

$$\mathbf{A}\mathbf{A}^+\mathbf{A} = \mathbf{A}$$

# $\mathbf{X}^T\mathbf{X}$ not invertible?

How do we get **w** such that

$$\mathbf{X}^T\mathbf{X}\mathbf{w} = \mathbf{X}^T\mathbf{y}$$

if $\mathbf{X}^T\mathbf{X}$ is not invertible?
Let

$$\mathbf{A} = \mathbf{X}^T\mathbf{X}$$

Let $\mathbf{A}^+$ be the *generalized inverse* of $\mathbf{A}$, i.e.:

$$\mathbf{A}\mathbf{A}^+\mathbf{A} = \mathbf{A}$$

### Proposition

If $\mathbf{A} = \mathbf{X}^T\mathbf{X}$ is not invertible, then $\hat{w} = \mathbf{A}^+\mathbf{X}^T\mathbf{y}$ is a solution to $\mathbf{X}^T\mathbf{X}\mathbf{w} = \mathbf{X}^T\mathbf{y}$.

# Computing the Generalized Inverse of **A**

Note $\mathbf{A} = \mathbf{X}^T\mathbf{X}$ is symmetric $\Rightarrow$ eigenvalue decomposition of $\mathbf{A}$:
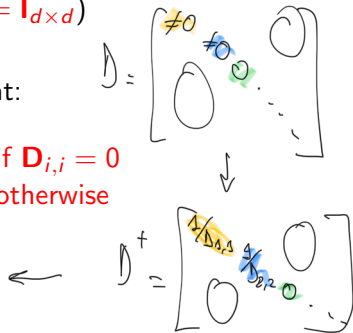
$$\mathbf{A} = \mathbf{V}\mathbf{D}\mathbf{V}^T$$

with

- **D**: diagonal matrix (entries = eigenvalues of **A**)
- **V**: orthonormal matrix ($\mathbf{V}^T\mathbf{V} = \mathbf{I}_{d \times d}$)

# Computing the Generalized Inverse of **A**

Note $\mathbf{A} = \mathbf{X}^T\mathbf{X}$ is symmetric $\Rightarrow$ eigenvalue decomposition of **A**:

$$\mathbf{A} = \mathbf{V}\mathbf{D}\mathbf{V}^T$$

with

- **D**: diagonal matrix (entries = eigenvalues of **A**)
- **V**: orthonormal matrix ($\mathbf{V}^T\mathbf{V} = \mathbf{I}_{d \times d}$)

Define $\mathbf{D}^+$ diagonal matrix such that:

$$\mathbf{D}_{i,i}^+ = \begin{cases} 0 & \text{if } \mathbf{D}_{i,i} = 0 \\ \frac{1}{\mathbf{D}_{i,i}} & \text{otherwise} \end{cases}$$

Let $\mathbf{A}^+ = \mathbf{V}\mathbf{D}^+\mathbf{V}^T$

Is it a generalized inverse for $A$?

Show $A A^+ A = A$

$$A A^+ A = \underbrace{V D V^T}_{A} \underbrace{V D^+ V^T}_{A^+} \underbrace{V D V^T}_{A}$$

since $V$ is orthonormal: $V^T V = I$

$$= V D D^+ D V^T$$

$$\underbrace{D D^+ D}_{D}$$

$$= V D V^T = A$$

11

Let $\mathbf{A}^+ = \mathbf{VD}^+\mathbf{V}^T$

Then

$$\begin{aligned}
\mathbf{AA}^+\mathbf{A} &= \mathbf{VDV}^T\mathbf{VD}^+\mathbf{V}^T\mathbf{VDV}^T \\
&= \mathbf{VDD}^+\mathbf{DV}^T \\
&= \mathbf{VDV}^T \\
&= \mathbf{A}
\end{aligned}$$

$\Rightarrow \mathbf{A}^+$ is a generalized inverse of $\mathbf{A}$.

Let $\mathbf{A}^+ = \mathbf{V}\mathbf{D}^+\mathbf{V}^T$

Then

$$\begin{aligned}
\mathbf{A}\mathbf{A}^+\mathbf{A} &= \mathbf{V}\mathbf{D}\mathbf{V}^T\mathbf{V}\mathbf{D}^+\mathbf{V}^T\mathbf{V}\mathbf{D}\mathbf{V}^T \\
&= \mathbf{V}\mathbf{D}\mathbf{D}^+\mathbf{D}\mathbf{V}^T \\
&= \mathbf{V}\mathbf{D}\mathbf{V}^T \\
&= \mathbf{A}
\end{aligned}$$

$\Rightarrow \mathbf{A}^+$ is a generalized inverse of $\mathbf{A}$.

**In practice**: the Moore-Penrose generalized inverse $\mathbf{A}^\dagger$ of $\mathbf{A}$ is used, since it can be efficiently computed from the Singular Value Decomposition of $\mathbf{A}$.

Your friend has developed a new machine learning algorithm for binary classification (i.e., $y \in \{-1, 1\}$) with 0-1 loss and tells you that it achieves a generalization error of only $0.05$. However, when you look at the learning problem he is working on, you find out that $\Pr_{\mathcal{D}}[y = 1] = 0.95$...

- Assume that $\Pr_{\mathcal{D}}[y = \ell] = p_\ell$. Derive the generalization error of the (dumb) hypothesis/model that *always* predicts $\ell$.

- Use the result above to decide if your friend's algorithm has learned something or not.

Assume we have the following training set $S$, where $\mathbf{x} = [x_1, x_2] \in \mathbb{R}^2$ and
$\mathcal{Y} = \{-1, 1\}$:
$S = \{([-3, 4], 1), ([2, -3], -1), ([-3, -4], -1), ([1, 1.5], 1)\}$.
Assume you decide to use $\mathcal{H} = \{h_1, h_2, h_3, h_4\}$ with
$h_1 = sign(-x_1 - x_2)$
$h_2 = sign(-x_1 + x_2)$
$h_3 = sign(x_1 - x_2)$
$h_4 = sign(x_1 + x_2)$
Your algorithm uses the ERM rule and the 0-1 loss.

- What model $h_S$ is produced in output by your ML algorithm?

- Assume the realizability assumption holds. What can you say about the
  generalization error $L_{\mathcal{D}}(h_S)$ of $h_S$?

Consider a linear regression problem, where $\mathcal{X} = \mathbb{R}^d$ and $\mathcal{Y} = \mathbb{R}$, with mean squared loss. The hypothesis set is the set of *constant* functions, that is $\mathcal{H} = \{h_a : a \in \mathbb{R}\}$, where $h_a(\mathbf{x}) = a$. Let $S = ((\mathbf{x}_1, y_1), \ldots, (\mathbf{x}_m, y_m))$ denote the training set.

- Derive the hypothesis $h \in \mathcal{H}$ that minimizes the training error.

- Use the result above to explain why, for a given hypothesis $\hat{h}$ from the set of all linear models, the coefficient of determination
  $R^2 = 1 - \frac{\sum_{i=1}^m (\hat{h}(\mathbf{x}_i) - y_i)^2}{\sum_{i=1}^m (y_i - \bar{y})^2}$ where $\bar{y}$ is the average of the $y_i, i = 1, \ldots, m$ is a
  measure of how well $\hat{h}$ performs (on the training set).

# Polynomial Models

Consider a regression problem.

*"linear in the parameters"*

Can we as hypothesis set the set of polynomials of degree $r$ with the tools we have already developed for linear regression?

Assume: $X = \mathbb{R}$

polynomial of degree $r$: $w_0 \cdot 1 + w_1 x + w_2 x^2 + \ldots + w_r x^r$

Given $x \in \mathbb{R}$, compute the following vector: (feature expansions)

$$\vec{x}' = \begin{bmatrix} 1 \\ x \\ x^2 \\ x^3 \\ \vdots \\ x^{r-1} \\ x^r \end{bmatrix} \Rightarrow \vec{w} = [w_0, w_1, \cdots, w_r]^T \Rightarrow \langle \vec{w}, \vec{x}' \rangle = w_0 \cdot 1 + w_1 x + w_2 x^2 + \ldots + w_r x^r$$

$\Rightarrow$ the hypothesis class of linear models for $\vec{x}'$ corresponds to the hypothesis class of polynomials of degree $r$ for $x$.

Given $\vec{x} \in \mathbb{R}^d$. $\vec{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix}$. You can use the following

feature expansion:

$\implies$ use linear models for $\vec{x}'$!

$$\vec{x}' = \begin{bmatrix} 1 \\ x_1 \\ x_1^2 \\ x_1^3 \\ \vdots^r \\ x_1^r \\ x_2 \\ x_2^2 \\ \vdots^r \\ x_2^r \\ \vdots \\ x_d \\ x_d^2 \\ \vdots^r \\ x_d^r \end{bmatrix}$$

Different feature expansion:

$\vec{x} \in \mathbb{R}^3$: $\vec{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$, $r^2$

$$\vec{x}' = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ x_3 \\ x_1^2 \\ x_2^2 \\ x_3^3 \\ x_1 \cdot x_2 \\ x_1 \cdot x_3 \\ x_2 \cdot x_3 \end{bmatrix}$$

$\implies$ build linear models for $\vec{x}'$

14