

Lecture 14

Linear Regression

- Sol HW 4
- HW 5 due
+ 24 h
- L VII Regression
+ Notes

fig-toy-linreg.png + figures
.....

Lecture Notes ~~VII~~ Regression

Marina Meilă
mmp@stat.washington.edu

Department of Statistics
University of Washington

~~May, 2023~~

Prediction problems ←

Linear regression ←

Linear regression for non-linear f

Reading: Ch.

Prediction

- ▶ **Data** $\mathcal{D} = \{(x^1, y^1), (x^2, y^2), \dots, (x^n, y^n)\}$
- ▶ **Inputs** $x^{1:n} \in \mathbb{R}$, or \mathbb{R}^d
- ▶ **Outputs** $y^{1:n}$

- ▶ **Goal** Learn/estimate $f(x)$ **predictor** for y

- ▶ By type of output
 - ▶ **Classification** if $y \in S_Y$ discrete
 - ▶ $y \in \{0, 1\}$ (or $\{\pm 1\}$) **binary classification**
 - ▶ $y \in \{1, 2, \dots, m\}$ **multiclass classification**

 - ▶ **Regression** if $y \in \mathbb{R}$ continuous

Part I Estimating a P on S $X \sim P$ on S
 + Model Selection

Part II Prediction $X \sim P_x$ on S_x
 $Y \sim P_{y|x}$ on S_y ← dependent on x

Prediction problems by S_y

- S_y discrete → classification
- $S_y \subset (-\infty, \infty)$ continuous → regression

Model of prediction

• $P_{y|x=x}$ uncertainty in y
 ↑ depends on input !!

• $X \sim P_x$ on S_x ← random x
 $Y \sim P_{y|x}$ on S_y ← random y given x uncertain
 ← ignoring / independently of it
 ← LEARN IT

$(x,y) \in S_x \times S_y$
 P_{xy} = joint
 P_x = marginal of X
 P_y = " of Y
 $P_{y|x}$ = conditional distribution of Y given x

Prediction Problem

▶ Data $D = \{(x^1, y^1), (x^2, y^2), \dots, (x^n, y^n)\}$

▶ Inputs $x^{1:n} \in \mathbb{R}$, or \mathbb{R}^d

▶ Outputs $y^{1:n}$ *output*

▶ Goal Learn/estimate $f(x)$ predictor for y

▶ By type of output *deterministic*

▶ **Classification** if $y \in S_Y$ discrete

▶ $y \in \{0, 1\}$ (or $\{\pm 1\}$) **binary classification**

▶ $y \in \{1, 2, \dots, m\}$ **multiclass classification**

▶ **Regression** if $y \in \mathbb{R}$ continuous

▶ **Model family** $\mathcal{F} = \{f_\theta : \mathbb{R}^d \rightarrow S_Y, \theta \in \Theta\}$

▶ $\mathcal{F} = \{ \text{linear functions} \}$ linear regression/classification

▶ $\mathcal{F} = \{ \text{polynomials of degree 2, 3, ...} \}$ polynomial regression/classification

▶ $\mathcal{F} = \{ \frac{1}{1+e^{-\beta^T x + \beta_0}} \}$ **logistic regression**

▶ **neural network** regression/classification

▶ **kernel** regression/classification

▶ $\mathcal{F} = \{ \text{monotonic functions} \}$ **isotonic regression**

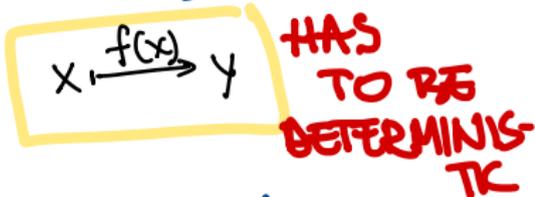
▶ **support vector** regression/classification

▶ regression/classification **trees** (and **random forests**)

TRAINING
(Learning)
(Estimation)

features
covariates
 $x = \text{input(s)}$
attributes

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix} = x \in \mathbb{R}^d$$



$$f : \mathbb{R}^d \rightarrow \mathbb{R}$$

parametric

non-param.

Linear regression (see also last 3 pages)

Given $\mathcal{D} = \{(x^i, y^i), i=1:n\}$

Model

- ▶ $y^i = \beta_0 + \beta_1 x^i + \epsilon^i$ for $x^{1:n} \in \mathbb{R}$ (univariate regression)
- ▶ $y^i = \beta_0 + \beta_1 x_1^i + \beta_2 x_2^i + \dots + \beta_d x_d^i + \epsilon^i$ for $x^{1:n} \in \mathbb{R}^d$ (multivariate regression)

$y \in \mathbb{R}$ is **output/response**

$x_j^{1:n}$ for $j = 1 : d$ are **input(s)/covariates/features/attributes/...**

$\beta_{1:d}$ are **regression coefficients**, β_0 is **intercept**

$\epsilon^{1:n} \in \mathbb{R}$ is **noise**, $\epsilon^{1:n} \sim N(0, \sigma^2)$ i.i.d.

Model

• $S_x = \mathbb{R}^d$

$S_y = \mathbb{R}$

sample spaces

$X \sim P_X$ on \mathbb{R}^d

$Y \sim P_{Y|X}$

deterministic

$Y = f_\beta(X) + \epsilon$

$N(0, \sigma^2)$
noise

* independent of X

** $f_\beta(x) = \beta_0 + \beta_1 x_1 + \dots + \beta_d x_d$

$= [\beta_0 \ \beta_1 \ \dots \ \beta_d] \cdot \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{bmatrix}$

*

dim.

want: parameters $\beta_{0:d}, \sigma^2$

← Estimate by ML.

Linear regression

Model

- ▶ $y^i = \beta_0 + \beta_1 x^i + \epsilon^i$ for $x^{1:n} \in \mathbb{R}$ (univariate regression)
 - ▶ $y^i = \beta_0 + \beta_1 x_1^i + \beta_2 x_2^i + \dots + \beta_d x_d^i + \epsilon^i$ for $x^{1:n} \in \mathbb{R}^d$ (multivariate regression)
- $y \in \mathbb{R}$ is **output/response**
 $x_j^{1:n}$ for $j = 1 : d$ are **input(s)/covariates/features/attributes/...**
 $\beta_{1:d}$ are **regression coefficients**, β_0 is **intercept**
 $\epsilon^{1:n} \in \mathbb{R}$ is **noise**, $\epsilon^{1:n} \sim N(0, \sigma^2)$ i.i.d.

deterministic



$$y = f_{\beta}(x) + \epsilon \sim N(f_{\beta}(x), \sigma^2) \equiv N(\mu_x, \sigma^2)$$

$N(0, \sigma^2)$
noise

- independent
of x

μ_x

deterministic
function of x

Model class $\mathcal{F} = \{ f_{\beta}, \beta_{0:d} \in \mathbb{R}^{d+1}, \sigma^2 > 0 \}$

linear functions of x

Solution by Maximum Likelihood

$$\mathcal{D} = \{ (x^i, y^i), i=1:n \}$$

► **Likelihood** Probability $y | X$, parameters

► **Parameters** $\beta_0, \beta_{1:d}, \sigma^2$

► $\beta_0^{ML}, \beta_{1:d}^{ML}, (\sigma^2)^{ML} = \underset{\beta_0, \beta_{1:d}, \sigma^2}{\operatorname{argmax}} l(y|X, \beta_0, \beta_{1:d}, \sigma^2)$

Max ^{log} Likelihood $(y^{1:n} | x^{1:n}, \beta_0, \beta_{1:d}, \sigma^2)$

want $P_{y|x}$
not also P_x

Before (Ch1):
 $l(\text{data} | \text{params})$

1. Likelihood

$$L(y^{1:n} | x^{1:n}, \beta_{0:d}, \sigma^2) =$$

$$= \prod_{i=1}^n e^{-\frac{(y^i - f_{\beta}(x^i))^2}{2\sigma^2}} \cdot \frac{1}{\sqrt{2\pi\sigma^2}}$$

$$y^i \sim N(f_{\beta}(x^i), \sigma^2)$$

$$1. \log-L$$

$$l(\beta, \sigma^2) = - \sum_{i=1}^n \left(\frac{(y^i - \beta^T x^i)^2}{2\sigma^2} \right) - \frac{1}{2} \cdot n \ln(2\pi\sigma^2)$$

$$2. \max_{\beta, \sigma^2} l(\beta, \sigma^2)$$

STAT \uparrow
CALCULUS \downarrow

2.

$$l = - \sum_{i=1}^n \frac{(y^i - \beta^T x^i)^2}{2\sigma^2} - \frac{1}{2} \cdot n \ln(2\pi\sigma^2)$$

max β, σ^2

$$\frac{\partial l}{\partial \beta_j} = -2 \cdot \frac{1}{2\sigma^2} \sum_{i=1}^n (y^i - \beta^T x^i) \cdot (-x_j^i) = 0$$

$j=0,1,\dots,d$

$$\sum_{i=1}^n y^i x_j^i = \beta^T \sum_{i=1}^n (x^i) x_j^i \text{ for all } j$$

Linear system
d+1 unknowns
d+1 eqn.

variables

$$\beta_0, \beta_1, \dots, \beta_d$$

$$f_{\beta}(x) = \beta_0 + \beta_1 x_1 + \dots + \beta_d x_d$$

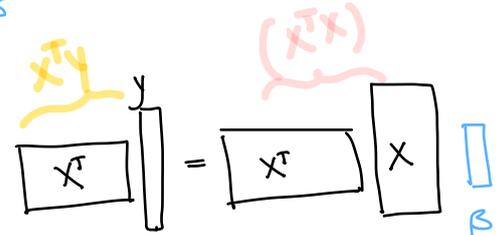
$$\frac{\partial}{\partial \beta_j} f_{\beta}(x) = x_j$$

$$X^T y = X^T X \beta$$

$$X = \begin{bmatrix} x_0^i & \dots & x_d^i \\ \dots & \dots & \dots \end{bmatrix}$$

n rows x (d+1) columns

$$y = \begin{bmatrix} y_1 \\ \dots \\ y_n \end{bmatrix}$$



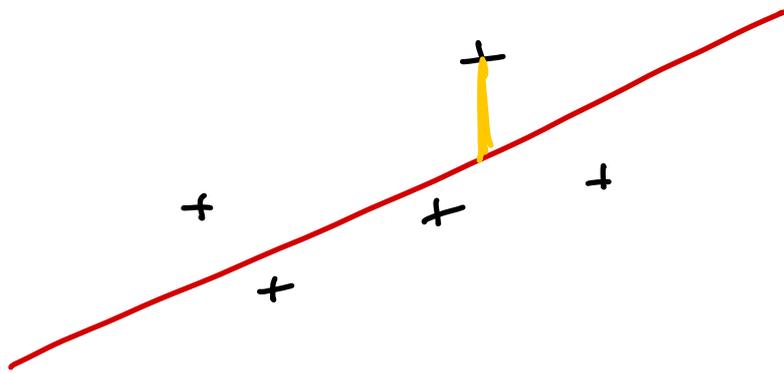
solve $\Rightarrow \beta^{ML} = (X^T X)^{-1} X^T y \in \mathbb{R}^{d+1}$

LEAST SQ SOLUTION

$$X^+ X = (X^T X)^{-1} (X^T X) = I$$

\hookrightarrow pseudo-inverse of X

X^+



$$f_b(x)$$

$$f_b(x^i) - y^i = \text{residual of } x^i$$

$$\min_{\beta} \sum_i (r_i)^2$$

More detailed derivation of Regression formula

$$\sum_i (y_i - \beta^T x^i) \cdot (-x_j^i) = 0 \quad j = 0:d$$

$$\sum_i \beta^T x^i x_j^i = \sum_i y_i x_j^i$$

$$\beta^T \sum_i \left(\begin{bmatrix} 1 \\ x^i \end{bmatrix} \begin{bmatrix} 1 & \dots & 1 \end{bmatrix} \right)$$

$(x^i)^T$

$$n \begin{bmatrix} 1 & \dots & 1 \\ x \end{bmatrix} \begin{bmatrix} x^i \\ y \end{bmatrix}$$

$y^T X \in \mathbb{R}^{d+1}$ row vector

$$\beta^T \left(\sum_i \begin{bmatrix} 1 \\ x^i \end{bmatrix} \begin{bmatrix} 1 & \dots & 1 \end{bmatrix} \right) \rightarrow X^T X$$

$$(X^T X) \beta = X^T y$$

$(d+1) \times (d+1)$

Linear regression model in matrix form

For a single data point (x^i, y^i)

$$\beta = \begin{bmatrix} \beta_1 \\ \beta_2 \\ \dots \\ \beta_d \end{bmatrix} \in \mathbb{R}^d \quad x^i = \begin{bmatrix} x_1^i \\ x_2^i \\ \dots \\ x_d^i \end{bmatrix} \in \mathbb{R}^d \quad (1)$$

Then,

$$y^i = \beta_0 + (x^i)^T \beta + \epsilon^i \quad (2)$$

For all \mathcal{D}

$$y = \begin{bmatrix} y^1 \\ y^2 \\ \dots \\ y^n \end{bmatrix} \in \mathbb{R}^n \quad X = \begin{bmatrix} (x^1)^T \\ (x^2)^T \\ \dots \\ (x^n)^T \end{bmatrix} \in \mathbb{R}^{n \times d} \quad \epsilon = \begin{bmatrix} \epsilon^1 \\ \epsilon^2 \\ \dots \\ \epsilon^n \end{bmatrix} \quad (3)$$

$$y = \beta_0 \mathbf{1} + X\beta + \epsilon \quad \text{with } \text{Cov}(\epsilon) = \sigma^2 I_d \quad (4)$$

The (log)-likelihood

- ▶ What is random? **the noise** $\epsilon^{1:n}$
- ▶ Express noise as function of $(x^{1:n}, y^{1:n})$

$$\epsilon^i = y^i - \beta_0 - \beta^T x^i \sim N(0, \sigma^2) \quad (5)$$

- ▶ Likelihood

- ▶ Let $p_{0, \sigma^2}(\epsilon) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{\epsilon^2}{2\sigma^2}} = N(\epsilon; 0, \sigma^2)$

- ▶ Then

$$L(\beta_0, \beta_{1:d}, \sigma^2) = \prod_{i=1}^n p_{0, \sigma^2}(\epsilon^i) \quad (6)$$

$$= \prod_{i=1}^n \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(\epsilon^i)^2}{2\sigma^2}} \quad (7)$$

$$= \prod_{i=1}^n \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(y^i - \beta_0 - \beta^T x^i)^2}{2\sigma^2}} \quad (8)$$

- ▶ **log-likelihood**

$$l(\beta_0, \beta_{1:d}, \sigma^2) = \quad (9)$$

$$= \sum_{i=1}^n \left\{ -\frac{1}{2} \ln \sigma^2 - \frac{1}{2} \ln(2\pi) - \frac{1}{2} (y^i - \beta_0 - \beta^T x^i)^2 \frac{1}{2\sigma^2} \right\} \quad (10)$$

$$= -\frac{n}{2} \ln \sigma^2 - \frac{1}{2} \ln(2\pi) + \text{constant} - \frac{1}{2\sigma^2} \sum_{i=1}^n (y^i - \beta_0 - \beta^T x^i)^2 \quad (11)$$

Maximizing the log-likelihood w.r.t β

- ▶ For simplicity, let $\beta_0 = 0$; hence $y^i = \beta^T x^i + \epsilon^i$
- ▶ log-likelihood

$$l(\beta_{1:d}, \sigma^2) = -\frac{n}{2} \ln \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (y^i - \beta_0 - \beta^T x^i)^2 + \text{constant} \quad (12)$$

- ▶ For any σ^2 ,

$$\underset{\beta}{\operatorname{argmax}} l(\sigma^2, \beta) = \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^n (y^i - \beta_0 - \beta^T x^i)^2 \quad (13)$$

a Least Squares Problem

- ▶ In matrix form $\min_{\beta} \|y - X\beta\|^2$
- ▶ Solution

$$\beta^{ML} = (X^T X)^{-1} X^T y \quad (14)$$

with $(X^T X)^{-1} X^T \equiv X^\dagger$ the **pseudoinverse** of X

- ▶ β^{ML} is **linear** in y !

vector calculus