

# 1. Probability, regression review

Friday, January 7, 2022 5:19 PM

## Probability Review

Definition/Convention A capital letter, like  $X$ , is usually a random variable

often written  $\underline{X}$  or  $\overline{X}$  by hand to make it clear

$$\mathbb{P}(X \geq .3) \text{ or } \text{Prob}(X \geq .3) \text{ or } P(X \geq .3)$$

mean the probability that  $X \geq .3$

A lowercase letter, like  $x$ , is usually a realization

$$\text{eg., } x = 0.3, \quad \mathbb{P}(X \geq x) = .5$$

Definition A time series  $X_t$  or  $(X_t)$  or  $X(t)$  is a sequence of random variables

If  $t \in N := \{0, 1, 2, 3, \dots\}$  or  $t \in \mathbb{Z} := \{0, \pm 1, \pm 2, \dots\}$

it's a discrete time series

If  $t \in \mathbb{R}$  or  $t \in [a, b] \subseteq \mathbb{R}$

it's a continuous time series.

our main focus  
for this class

A realization of a time series is  $X_t$  or  $(X_t)$  or  $X(t)$ , just a sequence.

We often call  $(X_t)$  a "time series" also

Definition Independence  $X_1, X_2, \dots, X_n$  are independent if

their cumulative distribution function (CDF) factors, namely

$$\underbrace{F(x_1, x_2, \dots, x_n)}_{\text{joint CDF}} = F_1(x_1) \cdot F_2(x_2) \cdot \dots \cdot F_n(x_n)$$

$$:= \mathbb{P}(X_1 \leq x_1, \dots, X_n \leq x_n)$$

Equivalently, the pdf/pmf factors

⚠ Independent variables are uncorrelated

but uncorrelated variables need not be independent \*

\*unless jointly normal.

Def Expectation,  $\mathbb{E}$ . Let  $X$  be a r.v.  
and Variance,  $\text{Var}$

- Continuous r.v. Let  $f_x$  be the pdf of  $X$ , then

$$\mathbb{E}[X] := \int x \cdot f_x(x) dx \quad \text{aka the mean, } \mu_x$$

↑ often written  
 $\mathbb{E}X$  or  $E[X]$  or  $EX$

- Discrete r.v. Let  $f_x$  be the pmf of  $X$ , then

$$\mathbb{E}[X] := \sum_x x \cdot f_x(x) = \mu_x$$

- Variance: in all cases,

$$\text{Var}[X] := \mathbb{E}[(X - \mu_x)^2] = \sigma_x^2$$

Note: not all distributions are cts or discrete (e.g., a mixture of the two).  
Similar def'n apply

- In all cases,

$$\textcircled{1} \quad \mathbb{E} \text{ is linear, } \mathbb{E}(aX + Y) = a \cdot \mathbb{E}(X) + \mathbb{E}(Y)$$

$$\textcircled{2} \quad \sigma_x^2 = \text{Var}[X] = \mathbb{E}[X^2] - \mathbb{E}[X]^2 \quad \text{very handy}$$

$\sigma^2 \geq 0$  via Jensen's Ineq.

$$\textcircled{3} \quad \mathbb{E}[X \cdot Y] = \mathbb{E}[X] \cdot \mathbb{E}[Y] \text{ only if } X, Y \text{ are uncorrelated}$$

(in particular, this is true if they're independent)

$$\textcircled{4} \quad \text{Var}[aX + b] = a^2 \cdot \text{Var}[X] \quad (\text{not linear})$$

$$\text{Var}[X + Y] = \text{Var}[X] + \text{Var}[Y] \text{ only if } X, Y \text{ are uncorrelated}$$

$$\textcircled{5} \quad \sigma = \sqrt{\text{Var}[X]} \text{ is the "standard deviation"}$$

Def Covariance

$$\text{Cov}(X, Y) := \mathbb{E}[(X - \mu_x)(Y - \mu_y)] = \mathbb{E}[X \cdot Y] - \underbrace{\mathbb{E}[X] \cdot \mathbb{E}[Y]}_{\mu_x \mu_y}$$

↑ over joint distribution

i.e.  $X, Y$  uncorrelated iff  $\text{Cov}(X, Y) = 0$

Properties

$$\textcircled{1} \quad \text{Cov}(X, X) = \text{Var}(X)$$

$$\textcircled{2} \quad \text{Cov}(aX + b, cY + d) = a \cdot c \cdot \text{Cov}(X, Y)$$

$$\textcircled{3} \quad \text{Cov}(X+Y, Z) = \text{Cov}(X, Z) + \text{Cov}(Y, Z)$$

Remark Intro probability classes focus on calculating  $E$ ,  $\text{Var}$ ,  $\text{Cov}$ ... using pdf/pmf  $f_x$  and/or cdf  $F_x$ .

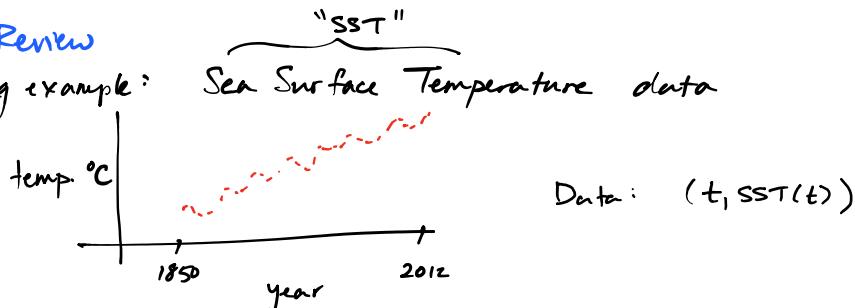
In more advanced usage, try to avoid using  $f_x$  and  $F_x$  directly until you really have to. Exploit properties (e.g., linearity of  $E$ )

### More prob. review

See Appendix A in Brockwell & Davis 3rd ed.

### Regression Review

Working example: Sea Surface Temperature data



We'll do ordinary least squares (OLS) regression

means there is not expected to be error in the independent variable (in this case, "year")

First, choose the independent variables aka covariates aka regressor

- The obvious choice is "year",  $t$ . What else?

We'll require our model to be linear in these covariates, so nonlinearity has to be explicitly added.

So could include  $t^2$  or  $\sin(t)$ , for example

- Let's stick with  $t$  and an offset (a constant).

so model: technically it's an "affine" function since we include the offset

$SST(t) = \text{linear function of time} + \text{error}$

$$= \beta_0 + \beta_1 t + \varepsilon(t)$$

$\theta$  is another common notation for  $\beta$

X as independent, Y as dependent is a universal convention.

$$\begin{array}{c} \text{---} \\ | \cdot \cdot \cdot \cdot \cdot \cdot \end{array} = \begin{array}{c} \text{---} \\ | \cdot \cdot \cdot \cdot \cdot \cdot \end{array} + \begin{array}{c} \text{---} \\ | \cdot \cdot \cdot \cdot \cdot \cdot \end{array}$$

to find the coefficients  $\vec{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix}$  we write

our model in vector form

163 years

$$\left\{ \begin{array}{l} \left[ \begin{array}{c} SST(1850) \\ \vdots \\ SST(2012) \end{array} \right] = \left[ \begin{array}{ccc} 1 & 1850 & \\ 1 & 1851 & \\ \vdots & \vdots & \\ 1 & 2012 & \end{array} \right] \left[ \begin{array}{c} \beta_0 \\ \beta_1 \end{array} \right] + \left[ \begin{array}{c} \varepsilon(1850) \\ \vdots \\ \varepsilon(2012) \end{array} \right] \end{array} \right.$$

matrix multiplication

$$(+) \quad \vec{Y} = X \cdot \vec{\beta} + \vec{\varepsilon} \leftarrow \text{"residual"}$$

now we're using the convention that capitals mean matrices, not random vectors

**OLS** in this case makes sense if we have the following assumptions:

(1) (+) is the true model, and

(2)  $\varepsilon(1850), \varepsilon(1851), \dots, \varepsilon(2012)$  are mutually independent

(3)  $\varepsilon(t) \sim N(0, \sigma^2)$  normally distributed  
(and "centered" aka mean 0)

then **OLS** is the maximum likelihood estimate (MLE)

for  $\vec{\beta}$ :

common convention:  
a hat  $\hat{\cdot}$  is used  
for estimators (something  
based on data)

$$\begin{aligned} \hat{\vec{\beta}}_{OLS} &= \underset{\vec{\beta}}{\operatorname{argmin}} \underbrace{\|\vec{Y} - X \cdot \vec{\beta}\|_2^2}_{=} \\ &= (\vec{Y} - X \vec{\beta})^\top (\vec{Y} - X \vec{\beta}) \\ &= \sum_{t=1850}^{2012} (SST(t) - \beta_0 - \beta_1 t)^2 \end{aligned}$$

Solvable in closed form by solving the normal equations

$$X^\top \cdot X \cdot \vec{\beta} = X^\top \vec{y}$$

so if  $X$  has full column rank,  $\vec{\beta} = (X^\top X)^{-1} X^\top \vec{y}$

and we can estimate  $\sigma^2$  also:

$$\hat{\sigma}^2 = \frac{1}{n-2} \|\vec{Y} - X \hat{\vec{\beta}}_{OLS}\|_2^2$$

you should almost never  
do this formula on  
a computer

# observations      size of  $\vec{\beta}$       "RSS"  
Residual Sum of Squares

$$R^2 = 1 - \frac{RSS}{SS_{tot}}$$

$R^2 \approx 1$  means most variability has been explained

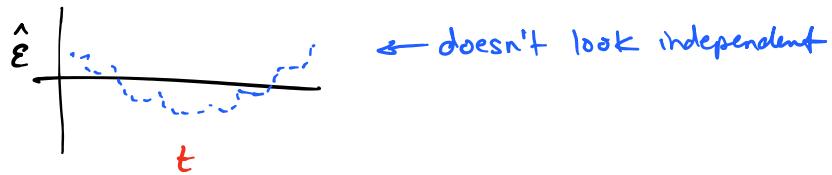
and define  $\hat{SST}(t) = \hat{\beta}_0 + \hat{\beta}_1 \cdot t$

and estimate residuals

$$\hat{\varepsilon}(t) = SST(t) - \hat{SST}(t)$$

Always a very good idea to plot residuals

we assumed true residuals were independent



We could model

$$SST(t) = \beta_0 + \beta_1 t + \underbrace{\beta_2 t^2}_{\text{new term}}$$

though if you add too many more terms, you're at risk of overfitting

In time series, we're going to discuss relaxing the assumption that  $\varepsilon(t)$  is independent

(both in a regression context and in other contexts)