

# Mathematical Statistics II

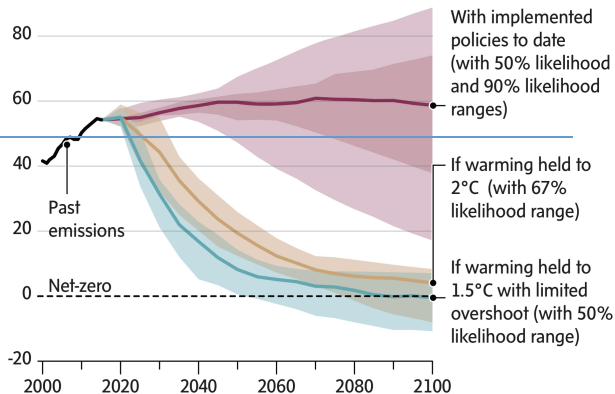
STA2212H S LEC0101

Week 11

March 28 2023

## Future emissions scenarios

Gigatonnes of CO<sub>2</sub> equivalent emissions per year

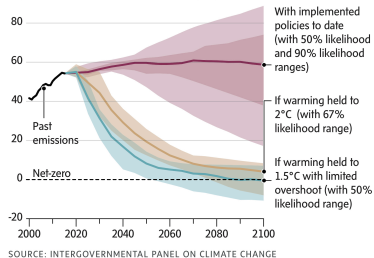


SOURCE: INTERGOVERNMENTAL PANEL ON CLIMATE CHANGE

# IPCC Report

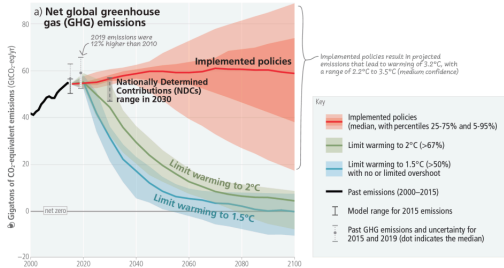
## Future emissions scenarios

Gigatonnes of CO<sub>2</sub> equivalent emissions per year



## Limiting warming to 1.5°C and 2°C involves rapid, deep and in most cases immediate greenhouse gas emission reductions

Net zero CO<sub>2</sub> and net zero GHG emissions can be achieved through strong reductions across all sectors



# Today

1. Recap
2. Directed acyclic graphs
3. Aspects of classification

## Upcoming

- April 4 10.00 – 13.00 Math Stat II Project presentations

please submit slides (pdf) by April 3

- April 11 9.30 – 12.30 Hydro Room 9016  
Informal discussion of large language models

# Recap

- observational studies can provide some evidence towards causality
- but care must be taken re confounding variables Simpson's "paradox"
- if all confounding variables are adjusted for, we have stronger evidence of the causal effect of a treatment on outcome
- this requires an assumption of "no unmeasured confounding"
- **Bradford-Hill guidelines** for strengthening support for causality  
in the absence of randomized treatment assignment

- one popular approach to causality is through the notion of **counterfactuals**
- the causal treatment effect is  $\theta = Y(1) - Y(0)$ ; the difference in outcome for an individual with  $X = 1$  compared to her outcome with  $X = 0$   $C_0, C_1$
- also called the causal risk difference
- since both outcomes cannot be observed, we must assume that in our data the units are “similar enough” that we can average over the treated and control to estimate  $\theta$
- $\alpha = E(Y \mid X = 1) - E(Y \mid X = 0)$  estimated by  $\bar{Y}_1 - \bar{Y}_0$  association

- causal risk difference  $\theta = Y(1) - Y(0)$
- causal risk difference as a function of an additional covariate  $Z$

$$\theta(z) = E(Y(1) \mid Z = z) - E(Y(0) \mid Z = z) = E(C_1 \mid Z = z) - E(C_0 \mid Z = z)$$

Thm 16.6: no unmeasured confounding

- causal regression function

continuous exposure  $X$

$$\theta(x) = E\{C(x)\}$$

- association function

$$r(x) = E(Y \mid X = x)$$

- Thm 16.4: If  $X$  assigned at random  $\theta(x) = r(x)$

- graphs can be useful for clarifying dependence relations among random variables

SM Markov random fields

- a **Directed Acyclic Graph** has random variables on the vertices and edges joining random variables

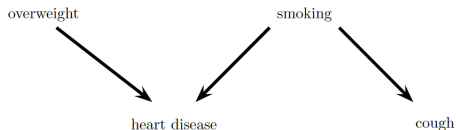


FIGURE 17.2. DAG for Example 17.4.

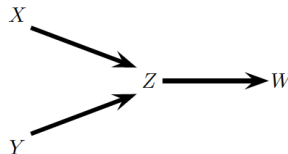


FIGURE 17.3. Another DAG.

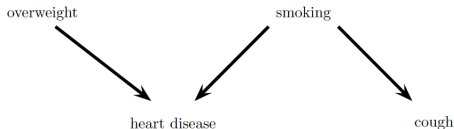


FIGURE 17.2. DAG for Example 17.4.

**17.4 Example.** Figure 17.2 shows a DAG with four variables. The probability function for this example factors as

$$\begin{aligned} & f(\text{overweight}, \text{smoking}, \text{heart disease}, \text{cough}) \\ &= f(\text{overweight}) \times f(\text{smoking}) \\ &\times f(\text{heart disease} \mid \text{overweight}, \text{smoking}) \\ &\times f(\text{cough} \mid \text{smoking}). \quad \blacksquare \end{aligned}$$

**17.5 Example.** For the DAG in Figure 17.3,  $\mathbb{P} \in M(\mathcal{G})$  if and only if its probability function  $f$  has the form

$$f(x, y, z, w) = f(x)f(y)f(z \mid x, y)f(w \mid z). \quad \blacksquare$$



- variables at parent nodes are potential causes for responses at child nodes
- probability distribution on a DAG represents causality if and only if the probability distribution is **Markov** wrt the DAG

AoS 17.5; HR Ch 6

- DAGs can be used to represent confounders

276 17. Directed Graphs and Conditional Independence

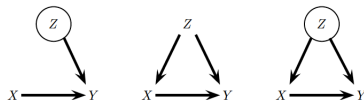


FIGURE 17.11. Randomized study; Observational study with measured confounders; Observational study with unmeasured confounders. The circled variables are unobserved.

- covariate space  $\mathcal{X} \subset \mathbb{R}^d$ ; prediction space  $\mathcal{Y} = \{0, 1\}$
- a **classification rule**

or  $\mathcal{Y} = \{1, \dots, K\}$

$$h : \mathcal{X} \rightarrow \mathcal{Y}$$

- data  $(X_1, Y_1), \dots, (X_n, Y_n)$ ;  $X_i \in \mathcal{X} \subset \mathbb{R}^d$ ,  $Y_i \in \mathcal{Y}$

supervised learning

350 22. Classification

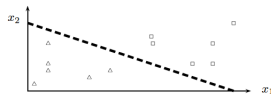


FIGURE 22.1. Two covariates and a linear decision boundary.  $\Delta$  means  $Y = 1$ .  $\square$  means  $Y = 0$ . These two groups are perfectly separated by the linear decision boundary; you probably won't see real data like this.

- loss function for classifier  $h$ :

$$L(h) = \text{pr}\{h(X) \neq Y\}$$

- empirical error rate

training error rate

$$\hat{L}_n(h) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{h(X_i) \neq Y_i\}$$

- special case  $\mathcal{Y} = \{0, 1\}$        $\hat{L}_n(h) =$

symmetric

- Bayes theorem

$$\text{pr}(Y = 1 \mid X = x) =$$

- loss function for classifier  $h$ :

$$L(h) = \text{pr}\{h(X) \neq Y\}$$

- empirical error rate

training error rate

$$\hat{L}_n(h) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{h(X_i) \neq Y_i\}$$

- special case  $\mathcal{Y} = \{0, 1\}$        $\hat{L}_n(h) =$

symmetric

- Bayes theorem

$$\text{pr}(Y = 1 \mid X = x) = \frac{f_1(x)\pi}{f_1(x)\pi + f_0(x)(1 - \pi)}$$

- Bayes theorem

$$\text{pr}(Y = 1 \mid X = x) = \frac{f_1(x)\pi}{f_1(x)\pi + f_0(x)(1 - \pi)} = r(x)$$

- Bayes classifier

$$h^*(x) = \begin{cases} 1 & r(x) > 1/2 \\ 0 & \text{otherwise} \end{cases} = \begin{cases} 1 & \pi f_1(x) > (1 - \pi)f_0(x) \\ 0 & \text{otherwise} \end{cases} =$$

- decision boundary

$$\mathcal{D} = \{x : \text{pr}(Y = 1 \mid X = x) = \text{pr}(Y = 0 \mid X = x)\}$$

- Thm 22.5: if  $h$  is another classification rule

$$L(h^*) \leq L(h)$$

- Bayes classifier:

$$h^*(x) = 1\{r(x) > 1/2\}, \quad r(x) = \frac{f_1(x)\pi}{f_1(x)\pi + f_0(x)(1-\pi)}$$

- multivariate normal

$$f_k(x) = \frac{1}{(2\pi)^{d/2}} |\Sigma_k|^{-1/2} \exp\left\{-\frac{1}{2}(x - \mu_k)^T \Sigma_k^{-1}(x - \mu_k)\right\}$$

- $r(x) > 1/2 \iff$

- $r(x) > 1/2 \iff$

$$(x - \mu_1)^T \Sigma_1^{-1} (x - \mu_1) < (x - \mu_0)^T \Sigma_0^{-1} (x - \mu_0) + \log(|\Sigma_0|/|\Sigma_1|) + 2 \log\{\pi/(1 - \pi)\}$$

- estimates of  $\pi_k, \mu_k, \Sigma_k$ :

- $\Sigma_0 = \Sigma_1$

- $\pi_0 = \pi_1$

- If  $\mathcal{Y} = \{1, \dots, K\}$ , the optimal classification rule is

$$\begin{aligned}h(x) &= \arg \max_k \text{pr}(Y = k \mid X = x) \\&= \arg \max_k \frac{\pi_k f_k(x)}{\sum_r \pi_r f_r(x)} \\&= \arg \max_k \pi_k f_k(x)\end{aligned}$$

- if  $f_k(x)$  is Gaussian, then

$$\begin{aligned}h^*(x) &= \arg \max_k \delta_k(x) \\ \delta_k(x) &= \frac{1}{2} \log |\Sigma_k| - \frac{1}{2} (x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k) + \log \pi_k\end{aligned}$$



- if  $\mathcal{Y} = \{0, 1\}$  and  $\Sigma_0 = \Sigma_1 = \Sigma$ , and  $\pi_1 = \pi_0 = 1/2$ , then

- if  $\mathcal{Y} = \{0, 1\}$  and  $\Sigma_0 = \Sigma_1 = \Sigma$ , and  $\pi_1 = \pi_0 = 1/2$ , then

- define  $w = S_W^{-1}(\bar{X}_1 - \bar{X}_0)$ ,  $S_W =$

- estimated Bayes classifier is

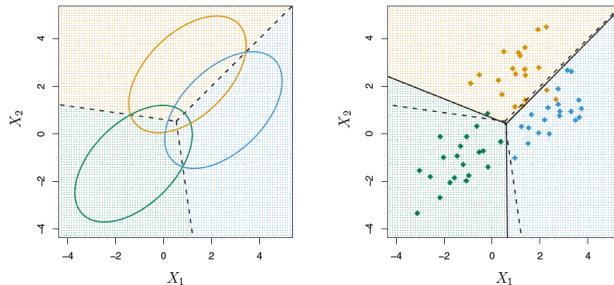
$$h^*(x) = \begin{cases} 0 & w^T x > m \\ 1 & w^T x < m \end{cases} \quad m = \frac{1}{2}(\bar{X}_0 + \bar{X}_1)$$

•

$$w = \arg \max_w \frac{w^T S_B w}{w^T S_W w}$$

- $w^T x \in \mathbb{R}$  maximizes between-group variation, relative to within-group variation

## 4.4 Linear Discriminant Analysis 143



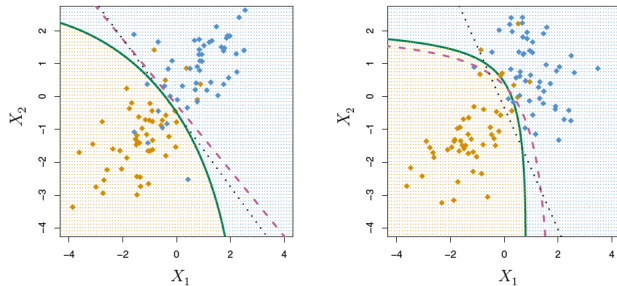
**FIGURE 4.6.** An example with three classes. The observations from each class are drawn from a multivariate Gaussian distribution with  $p = 2$ , with a class-specific mean vector and a common covariance matrix. Left: Ellipses that contain 95 % of the probability for each of the three classes are shown. The dashed lines are the Bayes decision boundaries. Right: 20 observations were generated from each class, and the corresponding LDA decision boundaries are indicated using solid black lines. The Bayes decision boundaries are once again shown as dashed lines.

- confusion matrix

	classified as 0	classified as 1
$y = 0$	277	25
$y = 1$	116	44

- misclassification rate  $(25 + 116)/(25 + 116 + 277 + 44) = 0.31$
- see ISLR §4.4 for discussion of changing cut-off from  $1/2$  to other thresholds
- changing the loss function to be asymmetric

## 150 4. Classification



**FIGURE 4.9.** Left: The Bayes (purple dashed), LDA (black dotted), and QDA (green solid) decision boundaries for a two-class problem with  $\Sigma_1 = \Sigma_2$ . The shading indicates the QDA decision rule. Since the Bayes decision boundary is linear, it is more accurately approximated by LDA than by QDA. Right: Details are as given in the left-hand panel, except that  $\Sigma_1 \neq \Sigma_2$ . Since the Bayes decision boundary is non-linear, it is more accurately approximated by QDA than by LDA.

- Model distribution of  $Y$ , given  $X$

- 

$$\text{pr}(Y_i = 1 \mid x_i) = \frac{\exp(\beta_0 + x_i^T \beta)}{1 + \exp(\beta_0 + x_i^T \beta)}$$

- Model distribution of  $Y$ , given  $X$

- 

$$\text{pr}(Y_i = 1 \mid \mathbf{x}_i) = \frac{\exp(\beta_0 + \mathbf{x}_i^T \boldsymbol{\beta})}{1 + \exp(\beta_0 + \mathbf{x}_i^T \boldsymbol{\beta})}$$

- equivalently

$$\log \left( \frac{\text{pr}(Y_i = 1 \mid \mathbf{x}_i)}{\text{pr}(Y_i = 0 \mid \mathbf{x}_i)} \right) = \beta_0 + \mathbf{x}_i^T \boldsymbol{\beta}$$

- compare LDA

$$\log \left( \frac{\text{pr}(Y_i = 1 \mid \mathbf{x}_i)}{\text{pr}(Y_i = 0 \mid \mathbf{x}_i)} \right) = \alpha_0 + \mathbf{x}_i^T \boldsymbol{\alpha}$$

$$\boldsymbol{\alpha}^T = (\mu_1 - \mu_0)^T \boldsymbol{\sigma}^{-1}$$

- $K$ -class again:

$$h(x) = \arg \max_k \Pr(Y = k \mid X = x) = \arg \max_k \pi_k f_k(x)$$

- we could estimate  $f_k(x)$  instead of assuming Gaussian, but  $X \in \mathbb{R}^d$
- pretend  $X_j$  are independent  $j = 1, \dots, d$
- one-dim density estimates from class  $k$ :

$$\hat{f}_k(x) = \prod_{j=1}^d \hat{f}_{kj}(x)$$

- class probabilities

$$\hat{\pi}_k = \frac{1}{n} \sum 1\{Y_i = k\}$$

- classifier

$$h(x) = \arg \max_k \hat{\pi}_k \hat{f}_k(x)$$

curse of  
dimensionality



- empirical error rate

misclassification rate

$$\hat{L}_n(h) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{h(X_i) \neq y_i\}$$

training set error

- increasing dimension of  $X_i$  decreases training error

e.g. adding variables in logistic regression

- true error rate

$$L(h) = \text{pr}\{h(X) \neq Y\}$$

- test error rate

$$L(\hat{h}) = \text{pr}\{\hat{h}(X_o) \neq Y_o \mid \mathcal{T}\}$$

- average test error rate

$$E_{\mathcal{T}}\{L(\hat{h})\}$$

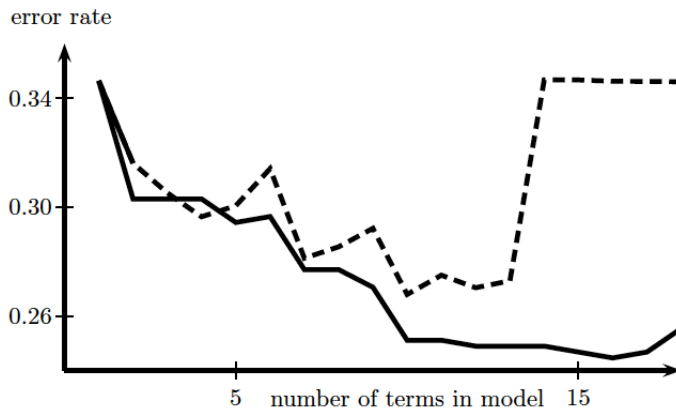


FIGURE 22.5. The solid line is the observed error rate and dashed line is the cross-validation estimate of true error rate.

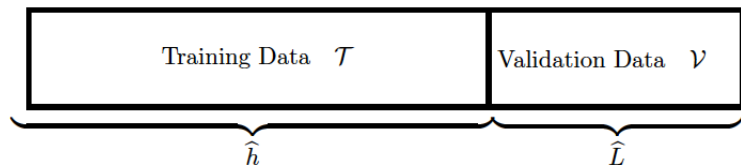


FIGURE 22.6. Cross-validation. The data are divided into two groups: the training data and the validation data. The training data are used to produce an estimated classifier  $\hat{h}$ . Then,  $\hat{h}$  is applied to the validation data to obtain an estimate  $\hat{L}$  of the error rate of  $\hat{h}$ .

- estimate test error rate with

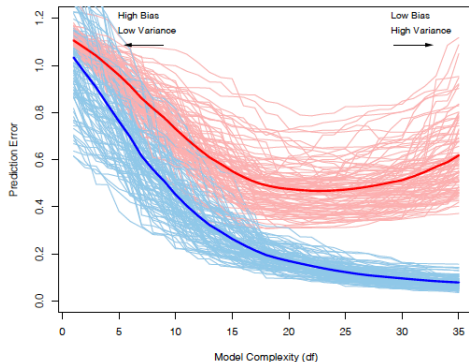
$$\hat{L}_m(h) \frac{1}{m} \sum_{i \in \mathcal{V}} \mathbf{1}\{\hat{h}(X_i) \neq Y_i\}$$

$K$ -fold cross-validation.

1. Randomly divide the data into  $K$  chunks of approximately equal size. A common choice is  $K = 10$ .
2. For  $k = 1$  to  $K$ , do the following:
  - (a) Delete chunk  $k$  from the data.
  - (b) Compute the classifier  $\hat{h}_{(k)}$  from the rest of the data.
  - (c) Use  $\hat{h}_{(k)}$  to predict the data in chunk  $k$ . Let  $\hat{L}_{(k)}$  denote the observed error rate.

3. Let

$$\hat{L}(h) = \frac{1}{K} \sum_{k=1}^K \hat{L}_{(k)}. \quad (22.33)$$

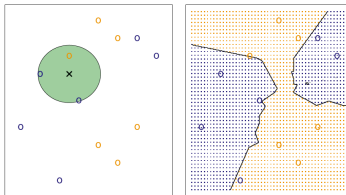


**FIGURE 7.1.** Behavior of test sample and training sample error as the model complexity is varied. The light blue curves show the training error  $\overline{\text{err}}$ , while the light red curves show the conditional test error  $\text{Err}_T$  for 100 training sets of size 50 each, as the model complexity is increased. The solid curves show the expected test error  $\text{Err}$  and the expected training error  $E[\overline{\text{err}}]$ .

- KNN:  $K$ – nearest neighbours
- choose a distance measure on  $\mathcal{X}$
- estimate  $\text{pr}(Y_1 | x)$  by averaging over “nearest neighbours” of  $x$

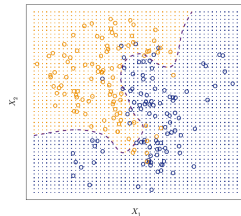
ISLR 2.2

40 2. Statistical Learning

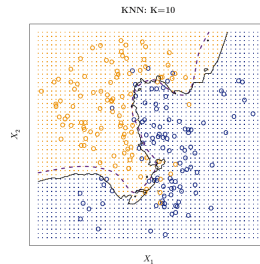


**FIGURE 2.14.** The KNN approach, using  $K = 3$ , is illustrated in a simple situation with six blue observations and six orange observations. Left: a test observation at which a predicted class label is desired is shown as a black cross. The three closest points to the test observation are identified, and it is predicted that the test observation belongs to the most commonly-occurring class, in this case blue. Right: The KNN decision boundary for this example is shown in black. The blue grid indicates the region in which a test observation will be assigned to the blue class, and the orange grid indicates the region in which it will be assigned to the orange class.

38 2. Statistical Learning



**FIGURE 2.13.** A simulated data set consisting of 100 observations in each of two groups, indicated in blue and in orange. The purple dashed line represents the Bayes decision boundary. The orange background grid indicates the region in which a test observation will be assigned to the orange class, and the blue background grid indicates the region in which a test observation will be assigned to the blue class.



**FIGURE 2.15.** The black curve indicates the KNN decision boundary on the data from Figure 2.13, using  $K = 10$ . The Bayes decision boundary is shown as a purple dashed line. The KNN and Bayes decision boundaries are very similar.

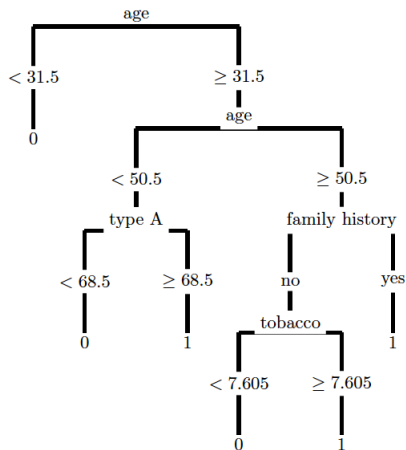


FIGURE 22.7. Smaller classification tree with size chosen by cross-validation.



- tree-based methods: bagging, boosting, random forests ISLR 8; ELSII 10
- support vector machines; kernelized SVMs ISLR 9; ESLII 12
- smoothing logistic regression ISLR 7.7.2
- neural networks ELSII 13
- double-descent in deep learning [link](#)