

# *Statistical Learning for Data Science*

*MSDS534*

## Lecture 1: Introduction and Review<sup>1</sup>

Department of Statistics  
Rutgers University

---

<sup>1</sup>Reproduction or redistribution without express written permission from Han Xiao is strictly prohibited.



# Acknowledgement

- Some of the figures in this presentation are taken from *An Introduction to Statistical Learning, with applications in R* (Springer, 2013) with permission from the authors: G. James, D. Witten, T. Hastie and R. Tibshirani.
- Some of the figures in this presentation are taken from *Elements of Statistical Learning* (Springer, 2009) with permission from the authors: T. Hastie, R. Tibshirani and J. Friedman.
- Some of the Lab codes in this presentation are taken from *An Introduction to Statistical Learning, with applications in R* (Springer, 2013) with permission from the authors: G. James, D. Witten, T. Hastie and R. Tibshirani.
- Some of these slides are from lectures by Prof. Zijian Guo, and Prof. Lingzhou Xue at Penn State.

## Data Sources

- UCI Machine Learning Repository:  
<https://archive.ics.uci.edu/ml/index.php>
- ESL: <https://web.stanford.edu/~hastie/ElemStatLearn/>
- ISL: <http://faculty.marshall.usc.edu/gareth-james/ISL/>

## Reading Assignments

- ESL: 1, §2.1–§2.6, §3.1–§3.6, §4.1 – §4.4.
- ISL: 1–4.

# Outline

1 AI, Machine Learning, and Statistical Learning

2 Review: Linear Regression

3 Review: Shrinkage methods

4 Review: Classification

# Dartmouth summer workshop on AI, 1956



Trenchard More, John McCarthy, Marvin Minsky, Oliver Selfridge, and Ray Solomonoff reunited at AI@50 conf.

# Legends

## Founding Fathers of AI

### **1956 Dartmouth Conference: The Founding Fathers of AI**



John McCarthy



Marvin Minsky



Claude Shannon



Ray Solomonoff



Alan Newell



Herbert Simon



Arthur Samuel

And three others...  
Oliver Selfridge  
(Pandemonium theory)  
Nathaniel Rochester  
(IBM, designed 701)  
Trenchard More  
(Natural Deduction)

- John McCarthy (1927-2011)
  - Coined the term "Artificial Intelligence" (1956)
  - ACM Turing Award (1971), National Medal of Science (1990)
- Marvin Minsky (1927-2016)
  - Co-founder of MIT CSAIL
  - ACM Turing Award (1969)
- Claude Shannon (1916-2001)
  - Father of Information Theory
  - National Medal of Science (1966), Shannon Award (1972)
- Ray Solomonoff (1926-2009)
  - Father of Algorithmic Probability & Algorithmic Information Theory
  - Kolmogorov Award (2003)
- Allen Newell (1927-1992)
  - ACM Turing Award (1975), National Medal of Science (1992)
- Herbert Simon (1916-2001)
  - ACM Turing Award (1975), National Medal of Science (1986)
  - Nobel Prize in Economics (1978)
- Arthur Samuel (1901-1990)
  - Coined the term "Machine Learning" (1959)
  - Worked with Donald Knuth on the TeX project

## ARTIFICIAL INTELLIGENCE

Early artificial intelligence stirs excitement.



## MACHINE LEARNING

Machine learning begins to flourish.



## DEEP LEARNING

Deep learning breakthroughs drive AI boom.



1950's

1960's

1970's

1980's

1990's

2000's

2010's

Since an early flush of optimism in the 1950s, smaller subsets of artificial intelligence – first machine learning, then deep learning, a subset of machine learning – have created ever larger disruptions.

- Artificial intelligence (coined by McCarthy)
  - The aspiration: human-imitative AI, the high-level or cognitive capability of humans to reason and to think.
  - Current development: engineering (pattern recognition, movement control...), statistics (finding pattern in data, making predictions, tests of hypothesis, decisions...).
  - Wiener coined **cybernetics**: a vision closely tied to operations research, statistics, pattern recognition, information theory, and control theory.
  - Wiener's agenda dominates under the banner of McCarthy's terminology.
- Machine learning
  - algorithmic field that blends ideas from statistics, computer science and many other disciplines
  - design algorithms that process data, make predictions, and help make decisions
- **Intelligent augmentation (IA)**. Computation and data are used to create services that augment human intelligence and creativity.
  - E.g. search engine, machine translation.
  - Will want computers to trigger new levels of human creativity.
- **Intelligent infrastructure (II)**. A web of computation, data, and physical entities that makes human environments more supportive, interesting, and safe.

<sup>2</sup>Michael I. Jordan. <https://hdsr.mitpress.mit.edu/pub/wot7mkc1/release/8>



## Hypothetical reasoning<sup>3</sup>

- People are smarter when they correctly use hypothetical reasoning. They access truths beyond those which can be inferred by simple direct measurement. They explore counterfactual situations they can imagine but not inhabit.
- Statistical hypothesis testing, the cornerstone of statistical inference, is a prime example of hypothetical reasoning.
- Formal statistical models offer another powerful mode of hypothetical reasoning with data.
- Two fields of special relevance.
  - **Causal inference.** Uses counterfactual reasoning to determine if the data we see might have been different if a certain hypothesized causal factor were silent.
  - **Robust inference.** Uses a hybrid of worst-case reasoning and statistical modeling to envision consequences of such hypothetical contamination and protect against them.
- True intelligence requires lots of (data-free) hypothetical reasoning about suspected causes, backed up by empirical checks.

---

<sup>3</sup>David Donoho. <https://hdsr.mitpress.mit.edu/pub/rim3pvdw/release/5>



# Role of Statistics<sup>4</sup>

- **Robustness.** Statistical reasoning and tools will be important.
  - Can we have “good enough” performance 99% of the time?
  - Can we be confident in our predictions?
  - How confident are our predictions?
- **Validity of algorithmic inference.** A deeper understanding of the data sources and the computations applied will be essential. Statistical quantification of uncertainty.
- **Fairness.**
  - Biased data collection yields biased results.
  - Statistical calculus of uncertainty, robustness, conditioning, (sub-)population quantities, and prediction errors are important.
- **Privacy.**
  - A sophisticated literature of algorithmic techniques under privacy constraints is growing.
  - More carefully integrated statistical reasoning is likely to yield tremendous benefits.

---

<sup>4</sup>Candès et al. <https://hdsr.mitpress.mit.edu/pub/djb16hz1/release/4>



# Outline

1 AI, Machine Learning, and Statistical Learning

2 Review: Linear Regression

3 Review: Shrinkage methods

4 Review: Classification

## Regression problem

- Quantitative generic output variable  $y$ .
- Generic input vector  $\mathbf{x} = (x_1, \dots, x_p)'$ .
- Use the regression function  $f(\mathbf{x})$  to predict  $y$ .
- Model the possible relationship between  $(\mathbf{x}, y)$  by a joint probability distribution.
- Take the loss function as

$$L[y, f(\mathbf{x})] = \mathbb{E}_{\mathbf{x}, y} [y - f(\mathbf{x})]^2,$$

then the optimal regression function is

$$f(\mathbf{x}) = \mathbb{E}_{y|\mathbf{x}}(y|\mathbf{x}).$$

- Typically we have a set of training data  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$ , from which we want to learn the regression function  $f(\mathbf{x})$ , or in statistical language, to estimate  $f(\mathbf{x})$ .

# Linear regression models

- Linear regression assumes that  $f(\mathbf{x})$  takes the form

$$f(\mathbf{x}) = \beta_0 + \sum_{j=1}^p x_j \beta_j.$$

- Terminology.

- $y$ : response, output, outcome, **dependent variable**
- $x_j$ : covariate, predictor, input, explanatory variable, feature, attribute, **independent variable**
- $\beta_j$ : parameter, regression coefficient

- The input  $x_j$  can come from different sources:

- quantitative inputs;
- transformations of quantitative inputs: log, square-root or square;
- basis expansions: polynomial, Fourier expansion, or splines;
- dummy variables derived from a discrete input;
- interactions, e.g.,  $x_3 = x_1 \cdot x_2$ .

- Linear means  $f$  is a linear function of  $\beta_j$ ,  $0 \leq j \leq p$ .

- $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2 + \beta_4 x_2^2$ .

## LSE: multiple linear regression

- Observations:  $(y_i, x_{i1}, \dots, x_{ip})$ ,  $1 \leq i \leq N$ .
- Define

$$\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix}, \quad \mathbf{X} = \begin{pmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1p} \\ 1 & x_{21} & x_{22} & \cdots & x_{2p} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{N1} & x_{N2} & \cdots & x_{Np} \end{pmatrix}, \quad \boldsymbol{\epsilon} = \begin{pmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_N \end{pmatrix}, \quad \boldsymbol{\beta} = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{pmatrix}$$

- Let  $\mathbf{x}_i = (1, x_{i1}, x_{i2}, \dots, x_{ip})'$  (**abuse of notation**), then for each individual observation

$$y_i = \mathbf{x}'_i \boldsymbol{\beta} + \epsilon_i, \quad 1 \leq i \leq N.$$

- The overall linear model can be represented as

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}.$$

- $\mathbf{X}$  is called **design matrix** or **model matrix**.
- Assume  $N > p + 1$  and  $\mathbf{X}$  has full rank  $p + 1$ .

## Least squares

- Minimize the sum of squared errors

$$\min_{\beta_0, \beta_1, \dots, \beta_p} \sum_{i=1}^N [y_i - (\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip})]^2.$$

- In matrix notation

$$\min_{\boldsymbol{\beta}} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2.$$

- Normal equation:

$$\mathbf{X}'\mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}'\mathbf{y}.$$

- Least squares estimator:

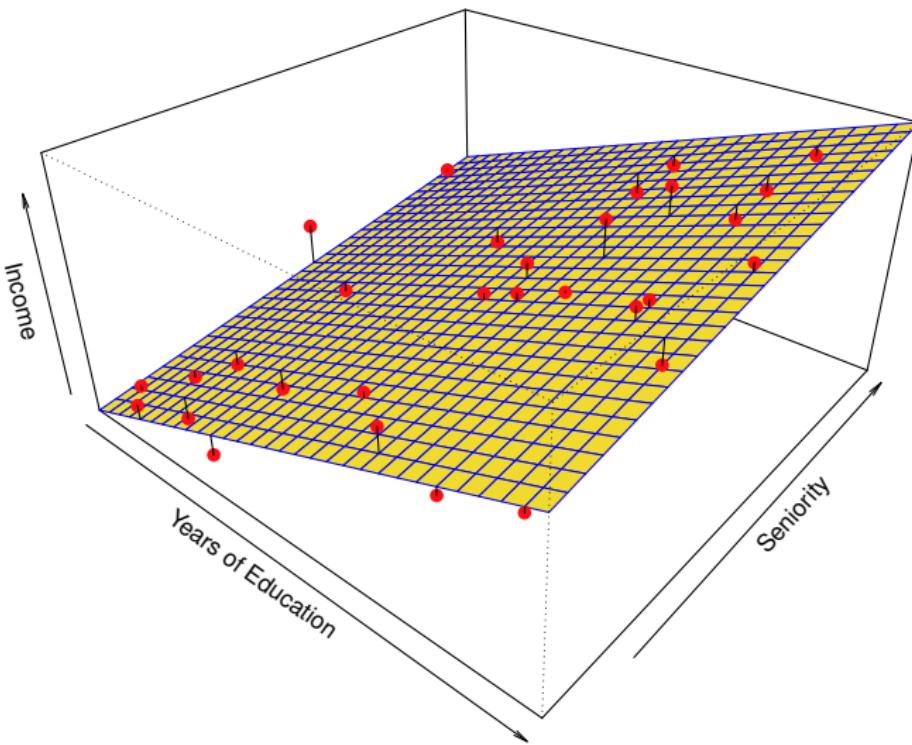
$$\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}.$$

- Hat matrix:  $\mathbf{H} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ .

- $\hat{y}_i = \mathbf{x}'_i \hat{\boldsymbol{\beta}}$ : fitted values. Vector form:  $\hat{\mathbf{y}} = \mathbf{H}\mathbf{y}$ .

- $\hat{\epsilon}_i = y_i - \hat{y}_i$ : residuals.

## Least squares: Income vs Education and Seniority<sup>5</sup>



<sup>5</sup>Figure 2.4 of ISL.

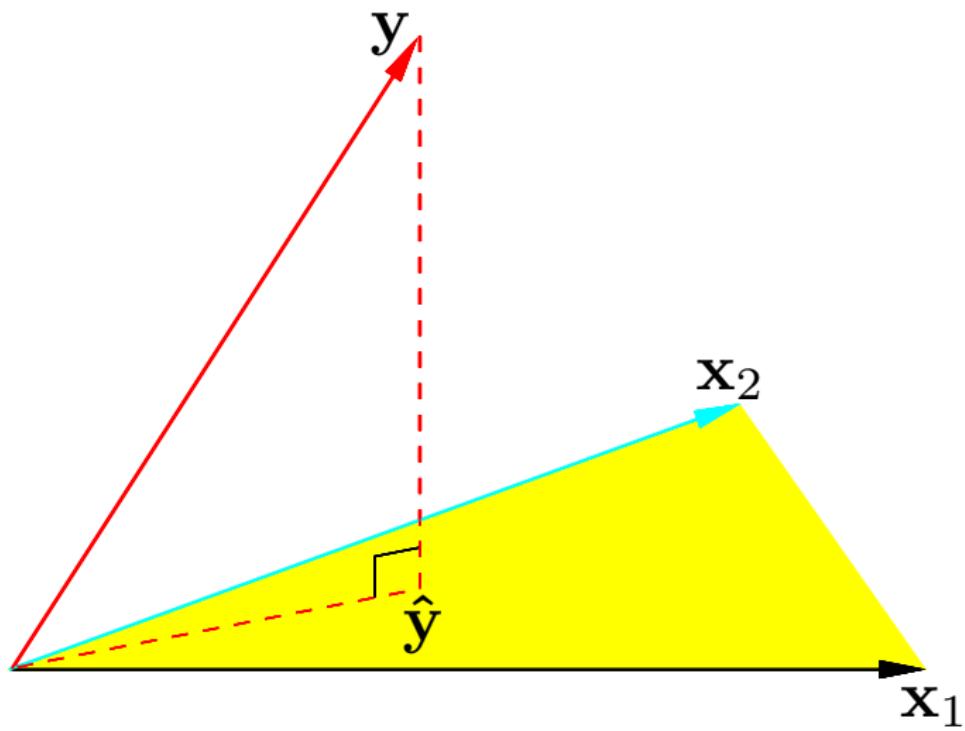
# Orthogonal projection

- Let  $\mathbf{x}^0, \mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^p$  be the  $p + 1$  column vectors of  $\mathbf{X}$ .
- The **column span** of  $\mathbf{X}$ , denoted by  $\text{col}(\mathbf{X})$  is the linear subspace consisting of all linear combinations of the form

$$b_0\mathbf{x}^0 + b_1\mathbf{x}^1 + b_2\mathbf{x}^2 + \cdots + b_p\mathbf{x}^p, \quad b_0, b_1, \dots, b_p \in \mathbb{R}.$$

- Least squares method finds a  $\hat{\boldsymbol{\beta}}$  such that  $\mathbf{X}\hat{\boldsymbol{\beta}}$  is the projection of  $\mathbf{y}$  onto the subspace  $\text{col}(\mathbf{X})$ .
- Let  $\mathbf{H}$  be the **projection matrix** onto  $\text{col}(\mathbf{X})$ .
  - Each column of  $\mathbf{H}$  belongs to  $\text{col}(\mathbf{X})$ , i.e.  $\mathbf{H} = \mathbf{X}\mathbf{C}'$  for some  $n \times (p + 1)$  matrix  $\mathbf{C}$ .
  - $\mathbf{H}\mathbf{X} = \mathbf{X}$ .
  - $\mathbf{H} = \mathbf{H}'$ , orthogonal projection.
- $\mathbf{H} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ .

## Orthogonal projection<sup>6</sup>



<sup>6</sup>Figure 3.2 of ESL.

## Variance estimator

- Recall the linear model is

$$y_i = \mathbf{x}'_i \boldsymbol{\beta} + \epsilon_i, \quad 1 \leq i \leq N.$$

- Assume  $\epsilon_1, \epsilon_2, \dots, \epsilon_N$  are independent and identically distributed (iid) with mean zero and variance  $\sigma^2$ .
- The estimator of  $\sigma^2$  is given by

$$\hat{\sigma}^2 = \frac{1}{N-p-1} \sum_{i=1}^N (y_i - \hat{y}_i)^2.$$

- Unbiased estimator, i.e.  $\mathbb{E}\hat{\sigma}^2 = \sigma^2$ .
- The covariance matrix of  $\hat{\boldsymbol{\beta}}$  is

$$\text{Cov}(\hat{\boldsymbol{\beta}}) = \sigma^2 (\mathbf{X}' \mathbf{X})^{-1}.$$

- It is estimated by  $\widehat{\text{Cov}}(\hat{\boldsymbol{\beta}}) = \hat{\sigma}^2 (\mathbf{X}' \mathbf{X})^{-1}$ .
- If the assumption is stronger:  $\epsilon_1, \epsilon_2, \dots, \epsilon_N$  are iid  $N(0, \sigma^2)$ , then  $\hat{\boldsymbol{\beta}}$  is the maximum likelihood estimator (MLE) of  $\boldsymbol{\beta}$ , and the MLE of  $\sigma^2$  is

$$\tilde{\sigma}^2 = \frac{1}{N} \sum_{i=1}^n (y_i - \hat{y}_i)^2.$$



# Outline

1 AI, Machine Learning, and Statistical Learning

2 Review: Linear Regression

3 Review: Shrinkage methods

4 Review: Classification

## Ridge regression

- The ridge regression solves the optimization problem

$$\hat{\beta}^{\text{ridge}} = \arg \min_{\beta} \left\{ \sum_{i=1}^N \left( y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 + \lambda \sum_{j=1}^p \beta_j^2 \right\}$$

for some  $\lambda \geq 0$ .

- An equivalent form is

$$\hat{\beta}^{\text{ridge}} = \arg \min_{\beta} \sum_{i=1}^N \left( y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2$$

subject to  $\sum_{j=1}^p \beta_j^2 \leq t$ ,

for some  $t \geq 0$ .

## Ridge regression

- Here  $\lambda \geq 0$  is a complexity parameter that controls the amount of shrinkage: the larger the value of  $\lambda$ , the greater the amount of shrinkage. The coefficients are shrunk toward zero (and each other).
- The ridge solutions are not equivariant under scaling of the inputs. Usually standardize the inputs before solving the optimization problem.
- The intercept  $\beta_0$  has been left out of the penalty term.
- Equivalent problem:
  - Center each  $x^j$ :  $x_{ij}^{(c)} := x_{ij} - \bar{x}_j$ .
  - Equivalent problem: estimate  $\beta_0$  by  $\bar{y} = \frac{1}{N} \sum_{i=1}^N y_i$ , and then solve

$$\min_{\beta_1, \dots, \beta_p} \left\{ \sum_{i=1}^N \left( y_i - \bar{y} - \sum_{j=1}^p x_{ij}^{(c)} \beta_j \right)^2 + \lambda \sum_{j=1}^p \beta_j^2 \right\}.$$

- By an abuse of notation, still use  $x_{ij}$  for  $x_{ij}^{(c)}$ .

# Ridge regression

- Pre-step.

- (i) The output vector  $\mathbf{y}$  is centered, that is,  $\sum_{i=1}^N y_i = 0$ ;
- (ii) Each predictor  $\mathbf{x}^j$ ,  $1 \leq j \leq p$  is normalized, i.e.

$$\sum_{i=1}^N x_{ij} = 0 \text{ and } \sum_{i=1}^N x_{ij}^2 = 1, \quad \forall 1 \leq j \leq p;$$

- (iii)  $\hat{\beta}_0 = \bar{y}$ .

- (iv) The input matrix  $\mathbf{X}$  has  $p$  (rather than  $p+1$ ) columns;

- Solve the problem (here  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)'$ )

$$\hat{\boldsymbol{\beta}}^{\text{ridge}} = \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^p} \left\{ \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda \|\boldsymbol{\beta}\|_2^2 \right\}.$$

- Ridge regression has a closed form solution

$$\hat{\boldsymbol{\beta}}^{\text{ridge}} = (\mathbf{X}'\mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}'\mathbf{y}.$$

- Compare with the LSE:  $\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y}$ .



# Lasso

- The LASSO solves the optimization problem

$$\hat{\beta}^{\text{lasso}} = \arg \min_{\beta} \left\{ \frac{1}{2} \sum_{i=1}^N \left( y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 + \lambda \sum_{j=1}^p |\beta_j| \right\}$$

for some  $\lambda \geq 0$ .

- An equivalent Lagrangian form is

$$\hat{\beta}^{\text{lasso}} = \arg \min_{\beta} \sum_{i=1}^N \left( y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2$$

subject to  $\sum_{j=1}^p |\beta_j| \leq t$ ,

for some  $t \geq 0$ .

## Solve lasso

- Pre-step.

- (i) The output vector  $\mathbf{y}$  is centered, that is,  $\sum_{i=1}^N y_i = 0$ ;
- (ii) Each predictor  $\mathbf{x}_{\cdot j}$ ,  $1 \leq j \leq p$  is normalized, i.e.

$$\sum_{i=1}^N x_{ij} = 0 \text{ and } \sum_{i=1}^N x_{ij}^2 = 1, \quad \forall 1 \leq j \leq p;$$

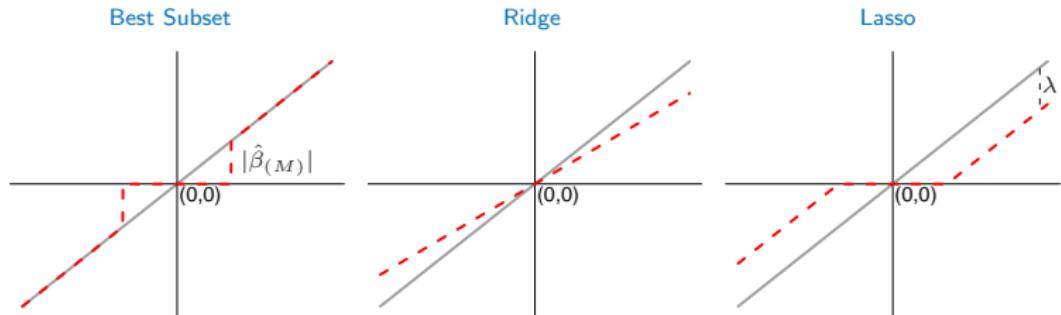
- (iii)  $\hat{\beta}_0 = \bar{y}$ .
- (iv) The input matrix  $\mathbf{X}$  has  $p$  (rather than  $p+1$ ) columns;
- Solve the problem (here  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)'$ )

$$\begin{aligned}\hat{\boldsymbol{\beta}}^{\text{lasso}} &= \arg \min_{\boldsymbol{\beta}} \left\{ \frac{1}{2} \sum_{i=1}^N \left( y_i - \sum_{j=1}^p x_{ij} \beta_j \right)^2 + \lambda \sum_{j=1}^p |\beta_j| \right\} \\ &= \arg \min_{\boldsymbol{\beta}} \left\{ \frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda \|\boldsymbol{\beta}\|_1 \right\}\end{aligned}$$

## Comparison: subset selection, ridge regression and lasso<sup>7</sup>

When the columns of  $\mathbf{X}$  are orthonormal, the formulas of different methods are given by

Estimator	Formula
Best subset (size $M$ )	$\hat{\beta}_j \cdot I( \hat{\beta}_j  \geq  \hat{\beta}_{(M)} )$
Ridge	$\hat{\beta}_j / (1 + \lambda)$
Lasso	$\text{sign}(\hat{\beta}_j)( \hat{\beta}_j  - \lambda)_+$



<sup>7</sup>Figure 3.11 of ESL.

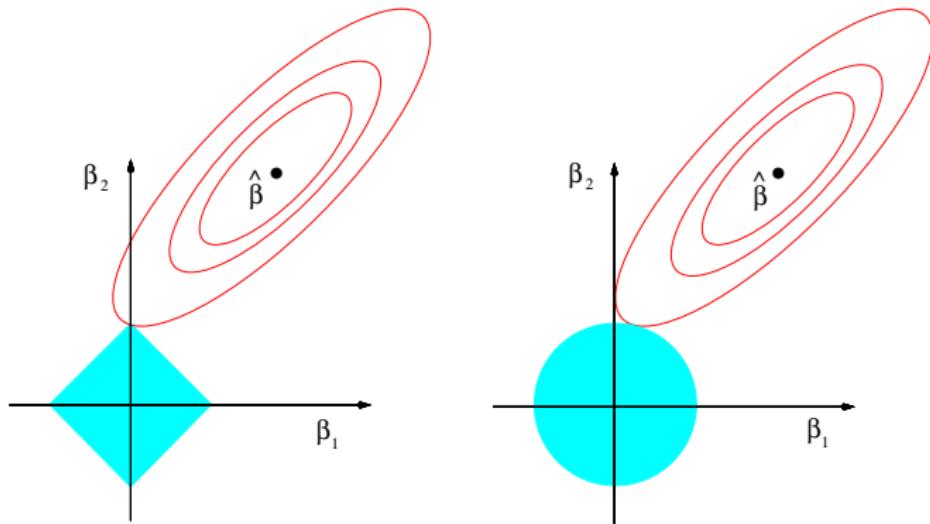
## Comparison: lasso and ridge regression<sup>8</sup>

$$\min_{\beta} \|\mathbf{y} - \mathbf{X}\beta\|_2^2$$

subject to

$$|\beta_1| + |\beta_2| \leq t$$

$$\sqrt{\beta_1^2 + \beta_2^2} \leq t$$



<sup>8</sup>Figure 3.11 of ESL.

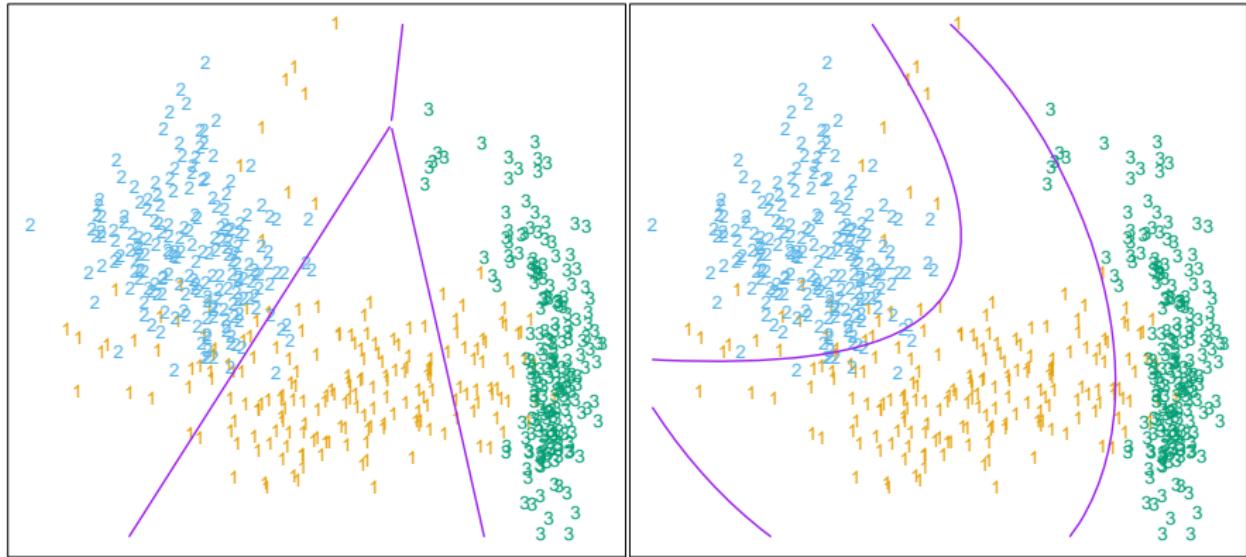
# Outline

- 1 AI, Machine Learning, and Statistical Learning
- 2 Review: Linear Regression
- 3 Review: Shrinkage methods
- 4 Review: Classification

## Classification: discriminative boundaries

- Each subject belongs to a class  $g$ , and comes with an **input**  $\mathbf{x}$ . Want to use  $\mathbf{x}$  to predict  $g$ .
- The **output**  $g$  takes value in a discrete set  $\mathcal{G}$  with cardinality  $|\mathcal{G}| = K$ , with each element representing a **class**. E.g.  $\{1, \dots, K\}$ ,  $\{0, 1, \dots, K - 1\}$ ,  $\{0, 1\}$ ,  $\{-1, 1\}$  etc.
- Divide the input space into a collection of regions labeled by classes. If we have an object  $(\mathbf{x}, g)$  whose **observed input**  $\mathbf{x}$  belong to the region labeled as  $k$ , then the **unobserved output**  $g$  is predicted as  $k$ .
- Associate each class  $k$  with a **discriminative function**  $\delta_k(\cdot)$ . For a subject with input  $\mathbf{x}$ ,  $\hat{g} = \arg \max_k \delta_k(\mathbf{x})$ .
- The **decision boundary** between class  $k$  and class  $l$  is the set of points such that  $\delta_k(\mathbf{x}) = \delta_l(\mathbf{x})$ .
- For many methods, either the discriminative function  $\delta_k(\mathbf{x})$  is linear in  $\mathbf{x}$ , or some monotone transformation of  $\delta_k(\mathbf{x})$  is linear in  $\mathbf{x}$ . As a result, the decision boundaries are linear in  $\mathbf{x}$ .

# Linear and nonlinear decision boundaries<sup>9</sup>



<sup>9</sup>Figure 4.1 of ESL.



## Discriminative method: logistic regression

- Consider the binary case, i.e.  $g \in \{0, 1\}$ ,  $\mathbf{x} \in \mathbb{R}^p$ . Logistic regression assumes:

$$\log \frac{P(g = 1|\mathbf{x})}{P(g = 0|\mathbf{x})} = \beta_0 + \boldsymbol{\beta}' \mathbf{x} \quad \text{i.e.}$$

- Equivalently, in terms of the discriminative functions:

$$\delta_0(\mathbf{x}) := P(g = 0|\mathbf{x}) = \frac{1}{1 + \exp(\beta_0 + \boldsymbol{\beta}' \mathbf{x})}$$
$$\delta_1(\mathbf{x}) := P(g = 1|\mathbf{x}) = \frac{\exp(\beta_0 + \boldsymbol{\beta}' \mathbf{x})}{1 + \exp(\beta_0 + \boldsymbol{\beta}' \mathbf{x})}.$$

- They are monotone functions of  $\beta_0 + \boldsymbol{\beta}' \mathbf{x}$ , and the decision boundary is

$$\{\mathbf{x} : \beta_0 + \boldsymbol{\beta}' \mathbf{x} = 0\}.$$

## Generative method: LDA

- For generative methods, the joint distribution of  $(\mathbf{x}, g)$  is of interest.
- Bayes Theorem.

$$P(g = k|\mathbf{x}) = \frac{P(\mathbf{x}|g = k)P(g = k)}{\sum_{l=1}^K P(\mathbf{x}|g = l)P(g = l)}.$$

- In general, assume the conditional density of  $\mathbf{x}$  given  $g = k$  is given by  $f_k(\mathbf{x})$ . Set  $\pi_k = P(g = k)$ . We have

$$P(g = k|\mathbf{x}) = \frac{\pi_k f_k(\mathbf{x})}{\sum_{l=1}^K \pi_l f_l(\mathbf{x})}.$$

- Assuming that  $f_k(\cdot)$  are multivariate normal distributions with different mean vectors, but SAME covariance matrix, then the decision boundary is linear in  $\mathbf{x}$ . Called **linear discriminant analysis (LDA)**.

## Generative method: linear discriminant analysis

- Gaussian discriminant analysis assumes

$$f_k(\mathbf{x}) = \frac{1}{(2\pi)^{N/2} \det(\Sigma_k)^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)' \Sigma_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) \right\}.$$

- Linear discriminant analysis also assumes  $\Sigma_k = \Sigma$  for all  $1 \leq k \leq K$ .
- To compare two classes  $k$  and  $l$ , it suffices to look at the log ratio

$$\log \frac{P(g = k | \mathbf{x})}{P(g = l | \mathbf{x})} = \log \frac{\pi_k}{\pi_l} - \frac{1}{2} (\boldsymbol{\mu}_k + \boldsymbol{\mu}_l)' \Sigma^{-1} (\boldsymbol{\mu}_k - \boldsymbol{\mu}_l) + \mathbf{x}' \Sigma^{-1} (\boldsymbol{\mu}_k - \boldsymbol{\mu}_l)$$

- The discriminant functions are given by

$$\delta_k(\mathbf{x}) = \log(\pi_k) - \frac{1}{2} \boldsymbol{\mu}_k' \Sigma^{-1} \boldsymbol{\mu}_k + \boldsymbol{\mu}_k' \Sigma^{-1} \mathbf{x}, \quad 1 \leq k \leq K.$$

- The parameters  $\{\pi_1, \mu_1, \dots, \pi_K, \mu_K, \Sigma\}$  are estimated by MLE.

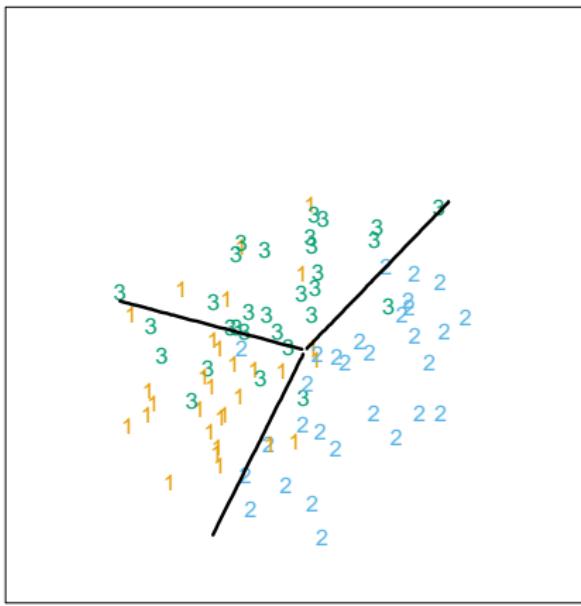
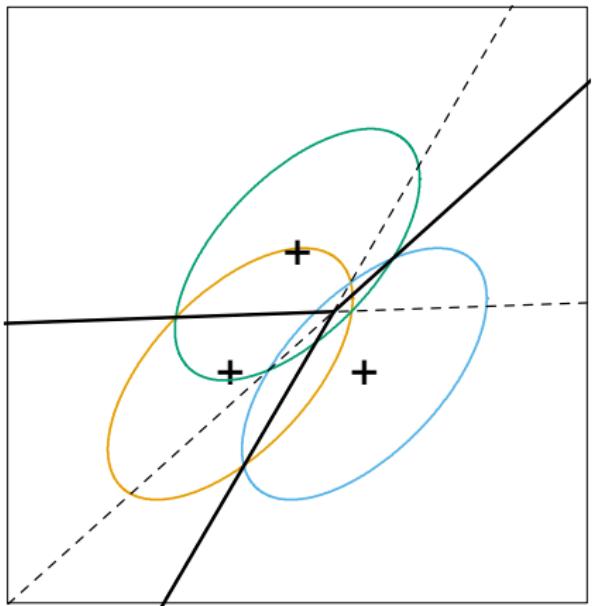
$$\hat{\pi}_k = \frac{1}{N} \sum_{i=1}^N I\{g_i = k\}$$

$$\hat{\boldsymbol{\mu}}_k = \frac{\sum_{i=1}^N \mathbf{x}_i I\{g_i = k\}}{\sum_{i=1}^N I\{g_i = k\}}$$

$$\hat{\Sigma} = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \hat{\boldsymbol{\mu}}_{g_i})^T (\mathbf{x}_i - \hat{\boldsymbol{\mu}}_{g_i}).$$



# LDA<sup>10</sup>



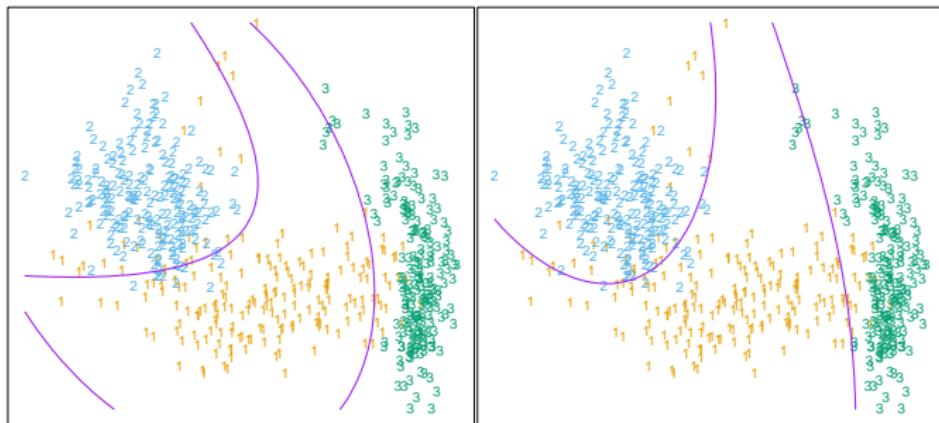
<sup>10</sup>Figure 4.5 of ESL.

## QDA and LDA in enlarged input space<sup>11</sup>

- If the  $\Sigma_k$  are not assumed to be equal, we need to use the **quadratic discriminant functions**

$$\delta_k(\mathbf{x}) = \log(\pi_k) - \frac{1}{2} \log[\det(\Sigma_k)] - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)' \Sigma_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k).$$

- Alternatively, we can perform LDA in the enlarged feature space:  
 $\{x_1, \dots, x_p, x_1^2, \dots, x_p^2, x_1 x_2, \dots, x_{p-1} x_p\}$ .



<sup>11</sup>Figure 4.6 of ESL.