

# Linear algebra, calculus, least squares & logistic regression

# Plan

2

- 25/04: Introduction
- 02/05: Linear algebra, calculus, least squares and logistic regression
- 09/05: SVM; k-fold cross-validation and boosting
- 16/05: CNNs; Backprop; Representation Learning; Regularisation; SGD
- 23/05: Image classification using Deep Learning models; Keras, Tensorflow and TF-tensorboard

**Please feel free to drop us an e-mail if you have any questions or suggestions.**

# Slides & Homework

You can find the presentation and homework exercises at:

<https://github.com/ink1/dl-training>

# References

- Lecture series on linear algebra:  
<https://www.khanacademy.org/math/linear-algebra>
- 3Blue1Brown's series on linear algebra  
(<https://youtu.be/kjBOesZCoqc>) and calculus  
(<https://youtu.be/WUvTyaaNkzM>)
- Part I of the Deep Learning book:  
<http://www.deeplearningbook.org/>
- Linear algebra explained in four pages:  
<https://minireference.com/blog/linear-algebra-tutorial/>
- Andrew Ng Machine Learning course at Stanford:  
<http://cs229.stanford.edu/syllabus.html>
- The Matrix Cookbook: <https://archive.org/details/imm3274>

# Linear Algebra

# Scalars

- A scalar is a single number: Integers, real numbers, rational numbers, ...
- Usually denoted with italic font:

*x, y, n*

# Vectors

- A vector is a 1-D array of scalars:

$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{pmatrix}$$

Usually denoted with lowercase, italic, bold font.

- A vector can contain real numbers, integers, ...
- Example notation of type and size:

$$\mathbf{x} = \begin{pmatrix} 2 \\ \pi \\ -\sqrt{2} \end{pmatrix} \in \mathbb{R}^3$$

# Matrices

- A matrix is a 2-D array of scalars:

The diagram shows a matrix  $X$  represented as a 2D array of scalars. The matrix is enclosed in large parentheses. The elements are arranged in rows and columns, with ellipses indicating continuation. A red oval highlights the first row, labeled "A row" in red text. A green oval highlights the  $m$ -th column, labeled "A column" in green text. A blue oval highlights the main diagonal, labeled "The main diagonal" in blue text. The matrix is labeled  $X =$  to its left.

$$X = \begin{pmatrix} x_{1,1} & x_{1,2} & \dots & x_{1,m} \\ x_{2,1} & x_{2,2} & \dots & x_{2,m} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n,1} & x_{n,2} & \dots & x_{n,m} \end{pmatrix}$$

- Usually denoted with uppercase, italic, bold font.
- $(X)_{i,j}$ : First index denotes the row, the second index the column, e.g.  $(X)_{4,2}$  denotes the scalar in the 4<sup>th</sup> row and 2<sup>nd</sup> column of matrix  $X$ .
- Example notation of type and shape:  $X \in \mathbb{R}^{n \times m}$



# Tensors

- A tensor is a an array of numbers of arbitrary dimensions.
- Scalars, vectors, and matrices are special cases of tensors.
  - A scalar is a zero-dimensional tensor
  - A vectors is a one-dimensional tensor
  - A matrix is a two-dimensional tensor

# Matrix Transpose

- The transpose can be thought of as a mirror image across the main diagonal.

$$(X^T)_{i,j} = X_{j,i}$$
$$X = \begin{pmatrix} x_{1,1} & x_{1,2} \\ x_{2,1} & x_{2,2} \\ x_{3,1} & x_{3,2} \end{pmatrix} \Rightarrow X^T = \begin{pmatrix} x_{1,1} & x_{2,1} & x_{3,1} \\ x_{1,2} & x_{2,2} & x_{3,2} \end{pmatrix}$$

$$((X)^T)^T = X$$

$$(X + Y)^T = X^T + Y^T$$

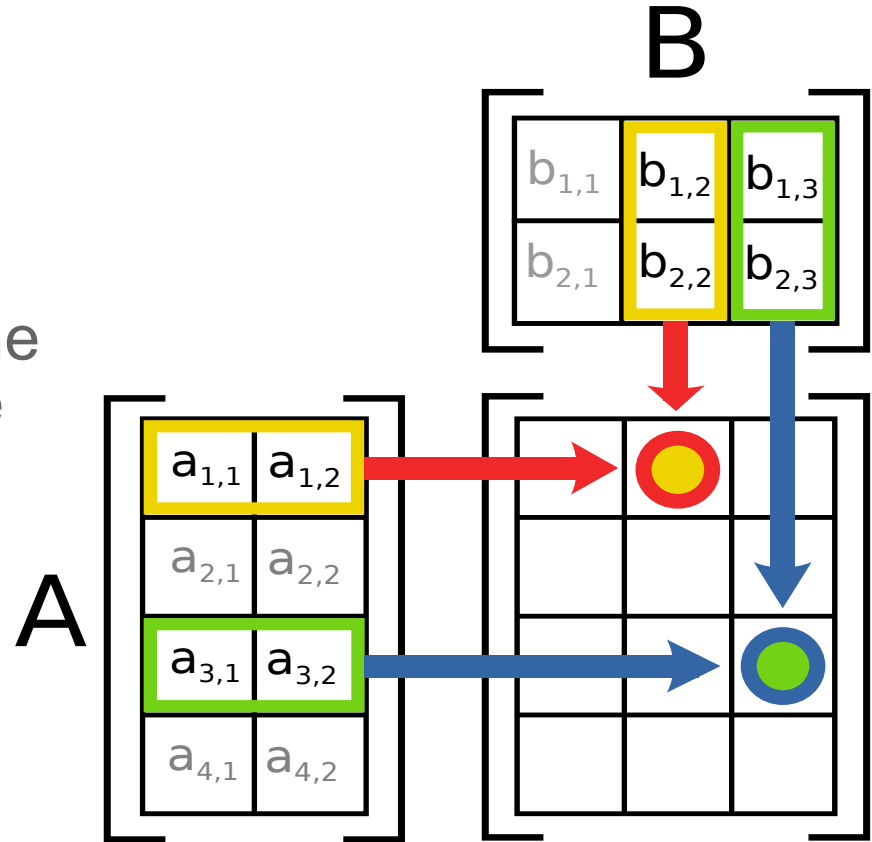
# Matrix Product

11

- The product of two matrices  $A \in \mathbb{R}^{n \times m}$  and  $B \in \mathbb{R}^{m \times p}$  is the matrix  $C = AB \in \mathbb{R}^{n \times p}$ , where

$$(C)_{i,j} = \sum_{k=1}^m (A)_{i,k} (B)_{k,j}$$

- Note that the number of columns in  $A$  and must equal the number of rows in  $B$ .



# Matrix product – Properties

12

- $(AB)^{\top} = B^{\top} A^{\top}$
- The matrix product is ...
  - distributive:  $A(B + C) = AB + AC$
  - associative:  $(AB)C = A(BC)$
  - is in general **not** commutative:  $AB \neq BA$

# Inner product – Vector-Vector Product

- Special case of matrix multiplication between a row-vector ( $1 \times n$ ) and a column-vector ( $n \times 1$ ).
- Given two vectors  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ , the inner product  $\mathbf{x}^\top \mathbf{y}$  is a scalar  $c \in \mathbb{R}$

$$\mathbf{x}^\top \mathbf{y} = \sum_{i=1}^n x_i y_i = c$$

- Also called **dot product**.

# Matrix-Vector product I

- The product of a matrix  $\mathbf{A} \in \mathbb{R}^{n \times m}$  and a vector  $\mathbf{x} \in \mathbb{R}^m$  is a vector  $\mathbf{y} = \mathbf{A}\mathbf{x} \in \mathbb{R}^n$ .
- If we express  $\mathbf{A}$  in terms of its rows, we can write the  $i$ -th entry of the product as the inner product of the  $i$ -th row of  $\mathbf{A}$  and  $\mathbf{x}$

$$\mathbf{y} = \begin{pmatrix} - & \mathbf{a}_1^\top & - \\ - & \mathbf{a}_2^\top & - \\ - & \vdots & - \\ - & \mathbf{a}_n^\top & - \end{pmatrix} \mathbf{x} = \begin{pmatrix} \mathbf{a}_1^\top \mathbf{x} \\ \mathbf{a}_2^\top \mathbf{x} \\ \vdots \\ \mathbf{a}_n^\top \mathbf{x} \end{pmatrix}$$

# Matrix-Vector product II

- We can also multiply a matrix on the left by a row vector.
- The product of a matrix  $\mathbf{A} \in \mathbb{R}^{n \times m}$  and a vector  $\mathbf{x} \in \mathbb{R}^m$  is a vector  $\mathbf{y}^\top = \mathbf{x}^\top \mathbf{A} \in \mathbb{R}^n$ .
- If we express  $\mathbf{A}$  in terms of its columns, we can write the  $i$ -th entry of the product as the inner product of the  $i$ -th column of  $\mathbf{A}$  and  $\mathbf{x}$

$$\mathbf{y}^\top = \mathbf{x}^\top \begin{pmatrix} | & | & \cdots & | \\ \mathbf{a}_1 & \mathbf{a}_2 & \cdots & \mathbf{a}_m \\ | & | & & | \end{pmatrix} = (\mathbf{x}^\top \mathbf{a}_1 \quad \mathbf{x}^\top \mathbf{a}_2 \quad \cdots \quad \mathbf{x}^\top \mathbf{a}_m)$$

# Matrix-Matrix product

- The product of two matrices  $\mathbf{A} \in \mathbb{R}^{n \times m}$  and  $\mathbf{B} \in \mathbb{R}^{m \times p}$  is the matrix  $\mathbf{C} = \mathbf{AB} \in \mathbb{R}^{n \times p}$ , where

$$(\mathbf{C})_{i,j} = \sum_{k=1}^m (\mathbf{A})_{i,k} (\mathbf{B})_{k,j}$$

- Thus, the entry  $(\mathbf{C})_{i,j}$  is the inner product of the  $i$ -th row of  $\mathbf{A}$  and the  $j$ -th column of  $\mathbf{B}$ .
- Alternatively, we can view matrix-matrix multiplication as a set of matrix-vector products by expressing  $\mathbf{B}$  by its columns

$$\mathbf{AB} = \mathbf{A} \begin{pmatrix} | & | & \cdots & | \\ \mathbf{b}_1 & \mathbf{b}_2 & \cdots & \mathbf{b}_p \\ | & | & & | \end{pmatrix} = \begin{pmatrix} | & | & \cdots & | \\ \mathbf{Ab}_1 & \mathbf{Ab}_2 & \cdots & \mathbf{Ab}_p \\ | & | & & | \end{pmatrix}$$



# Identity Matrix

- The **identity matrix**, denoted  $I_n \in \mathbb{R}^{n \times n}$ , is a square matrix with ones on the diagonal and zeros everywhere else

$$I_n = \begin{pmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & 0 & 0 \\ \vdots & & \ddots & \vdots \\ 0 & \cdots & 0 & 1 \end{pmatrix}$$

- It has the property that for any matrix  $A \in \mathbb{R}^{m \times n}$ ,

$$AI_n = A = I_n A$$

# Diagonal matrix

- A **diagonal matrix** is a matrix where all non-diagonal elements are 0. This is typically denoted

$$D = \text{diag}(d_1, d_2, \dots, d_n) = \begin{pmatrix} d_1 & 0 & \dots & 0 \\ 0 & d_2 & 0 & 0 \\ \vdots & & \ddots & \vdots \\ 0 & \dots & 0 & d_n \end{pmatrix}$$

# Inverse of a matrix

- The inverse of a square matrix  $A \in \mathbb{R}^{n \times n}$  is denoted  $A^{-1}$ , and is the unique matrix such that

$$A^{-1}A = I_n = AA^{-1}$$

- Matrix can't be inverted if ...
  - $A$  has more rows than columns,
  - $A$  has more columns than rows,
  - $A$  has redundant rows/columns (“linearly dependent”, “low rank”, “singular”).

# Inverse of a matrix – Properties

- The following properties assume that matrices  $A, B \in \mathbb{R}^{n \times n}$  are invertible.

$$(A^{-1})^{-1} = A$$

$$(AB)^{-1} = B^{-1} A^{-1}$$

$$(A^{-1})^T = (A^T)^{-1}$$

# Linear System of Equations I

- Consider the following system of equations with the unknowns  $x_1$  and  $x_2$ :

$$4x_1 - 5x_2 = -13$$

$$-2x_1 + 3x_2 = 9$$

- In matrix notation, we can write the system more compactly as  $A\mathbf{x} = \mathbf{b}$

$$\begin{pmatrix} 4 & -5 \\ -2 & 3 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} -13 \\ 9 \end{pmatrix}$$

# Linear System of Equations I

- Consider the following system of equations

$$Ax = b$$

$$\begin{pmatrix} 4 & -5 \\ -2 & 3 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} -13 \\ 9 \end{pmatrix}$$

- Multiply by  $A^{-1}$  on both sides to obtain

$$A^{-1}Ax = A^{-1}b$$

$$I_2x = A^{-1}b$$

$$x = A^{-1}b$$

# Linear System of Equations II

- A linear system of equations can have
  - No solution
  - Many solutions
  - Exactly one solution

# Vector Norm

- A norm of a vector  $\mathbf{x} \in \mathbb{R}^n$  denoted  $\|\mathbf{x}\| \in \mathbb{R}$  is a measure of the “length” of a vector.
- A norm must satisfy four properties:
  1. Non-negativity:  $\|\mathbf{x}\| \geq 0$
  2. Definiteness:  $\|\mathbf{x}\| = 0 \iff \mathbf{x} = \mathbf{0}$
  3. Homogeneity:  $\forall \alpha \in \mathbb{R}, \|\alpha \mathbf{x}\| = |\alpha| \|\mathbf{x}\|$
  4. Triangle inequality:  $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$



# Common Vector Norms

- $\ell_1$ -norm

$$\|\mathbf{x}\|_1 = \sum_{i=1}^n |x_i|$$

- $\ell_2$ -norm

$$\|\mathbf{x}\|_2 = \sqrt{\sum_{i=1}^n x_i^2}$$

- Max/Infinite-norm

$$\|\mathbf{x}\|_\infty = \max_i |x_i|$$

# Special Matrices and Vectors

- The vector  $\mathbf{x} \in \mathbb{R}^n$  is a **unit vector** if

$$\|\mathbf{x}\|_2 = 1$$

- The square matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$  is **symmetric** if

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

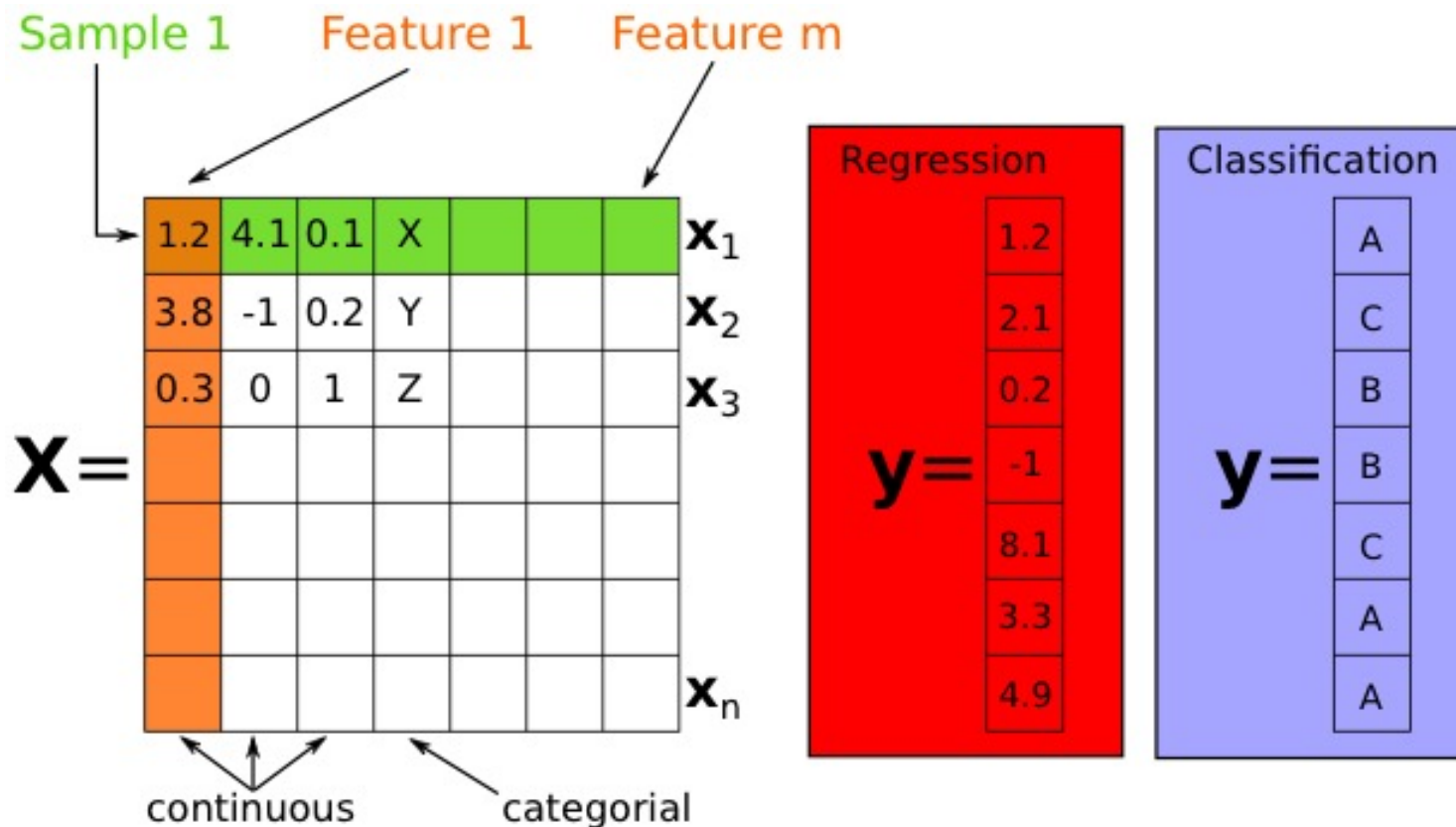
- The square matrix  $\mathbf{A} \in \mathbb{R}^{n \times n}$  is **orthogonal** if

$$\mathbf{A}\mathbf{A}^\top = \mathbf{A}^\top\mathbf{A} = \mathbf{I}_n$$

$$\mathbf{A}^{-1} = \mathbf{A}^\top$$

# Linear Models

# Training Data



- In machine learning, we usually have a matrix of measurements/features  $X \in \mathbb{R}^{n \times m}$  and a corresponding list of values  $y \in \mathbb{R}^n$  we want to predict.

# A Dataset

## Definitions

- A training **sample**  $x_i$  consists of  $m$  **features**  $(x_{i1}, \dots, x_{im})^T$  and is associated with **output**  $y_i$ .
- Each feature and the output can either be **continuous** (a number) or **discrete** (from a predefined set of values).
- If the output is **continuous**, we perform **regression** and if it is **discrete**, **classification**.
- The **training set**  $\mathcal{T} = (x_i, y_i)$  is comprised of  $n$  samples ( $i = 1, \dots, n$ ).
- Let  $X$  indicate a matrix where the  $i$ -th row corresponds to the  $i$ -th sample and  $\mathbf{y} = (y_1, \dots, y_n)^T$  the vector of all outputs.

# Body fat dataset

- Accurate measurement of body fat is inconvenient/costly and it is desirable to have easy methods of estimating body fat that are not inconvenient/costly.
- The dataset lists estimates of the percentage of body fat determined by underwater weighing and various body circumference measurements for 252 men.
- If we are able to accurately predict body fat from easy to obtain circumference measurements, we found an easy way to estimate body fat that is cost-effective and convenient.

# The Model

- Denote the  $i$ -th row of  $\mathbf{X} \in \mathbb{R}^{n \times m}$  as the  $i$ -th **feature vector**  $\mathbf{x}_i \in \mathbb{R}^m$ .
- By choosing a model, we determine how information captured by a feature vector  $\mathbf{x}_i$  is used to form a prediction  $\hat{y}_i \in \mathbb{R}$ .
- A machine learning **model**  $\mathcal{F}(\mathbf{x}; \Theta)$  maps a feature vector  $\mathbf{x}_i \in \mathbb{R}^m$  to a **prediction**

$$\hat{y}_i = \mathcal{F}(\mathbf{x}_i; \Theta) \in \mathbb{R}$$

- $\Theta$  is a set of **unknown parameters** we want to learn from data.

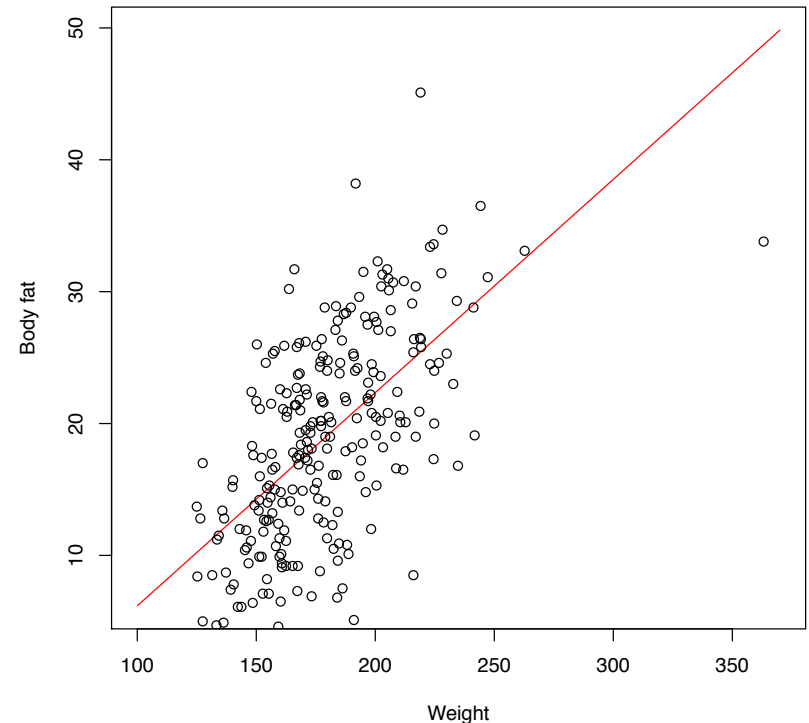
# Linear Model

## Definitions

### Definition

$$y_i = \beta_0 + \beta_1 x_{i,1} + \beta_2 x_{i,2} + \cdots + \beta_m x_{i,m} + \epsilon_i = \beta_0 + \mathbf{x}_i^\top \boldsymbol{\beta} + \epsilon$$

- The  $\boldsymbol{\beta}$  parameters are **coefficients** or weights of the features.
- $\boldsymbol{\beta}$  is to be **estimated** from training data.
- The errors  $\epsilon_i$  are independently and identically distributed (i.i.d.) with  $E(\epsilon_i) = 0$  and  $\text{Var}(\epsilon_i) = \sigma^2$ .

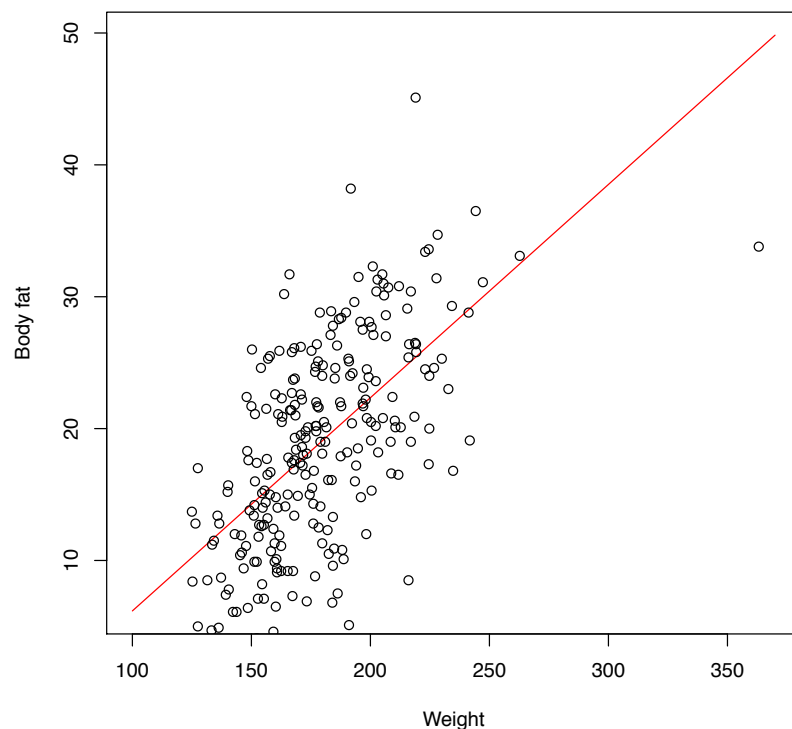




# Linear Model

## Coefficients

- Each **feature** is associated with one **coefficient**  $\beta_j$ .
- In addition, the coefficient  $\beta_0$  denotes the **intercept** (or bias).
- Estimates are denoted by a **hat**:  $\hat{\beta}_j$  denotes the estimate of the coefficient of the  $j$ -th feature.
- In the example to the right,  $\beta_0 = -9.995$  ( $y$ -intercept) and  $\beta_1 = 0.1617$  (slope; coefficient of *weight* feature).



# The Loss Function

- **Training data** consists of a matrix of features  $X \in \mathbb{R}^{n \times m}$  and a corresponding list of values  $y \in \mathbb{R}^n$  we want to predict.
- A **model**  $\mathcal{F}(x; \Theta)$  maps a feature vector  $x_i \in \mathbb{R}^m$  to a prediction  $\hat{y}_i = \mathcal{F}(x_i; \Theta) \in \mathbb{R}$  ( $\Theta$  are unknown parameters).
- We need to choose a function  $\mathcal{L}(y, \hat{y})$  that measures how well our model is approximating our training data.

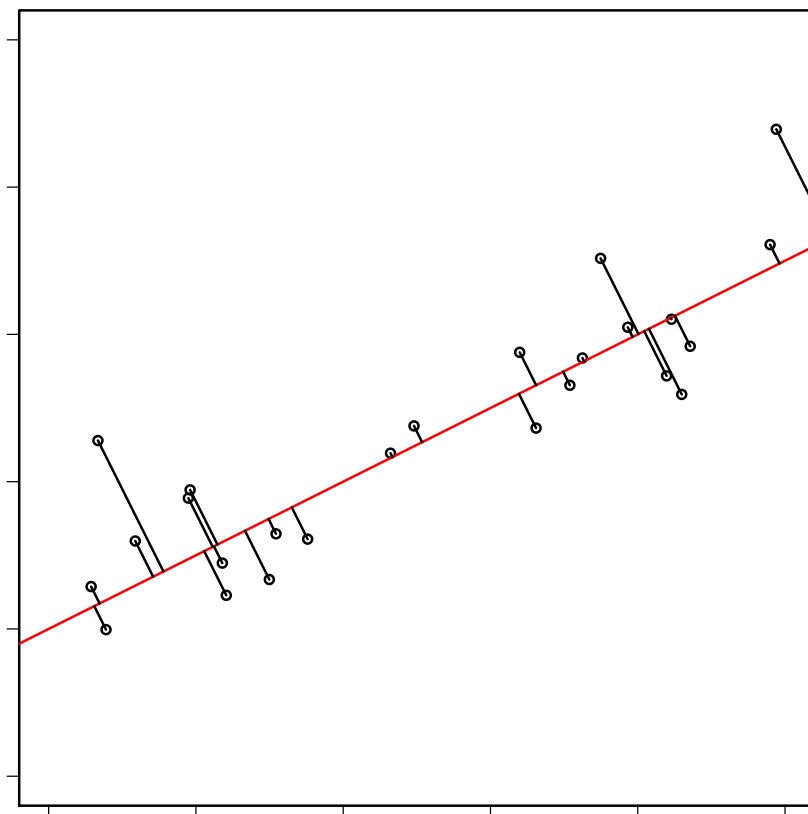
# Linear Model

## Loss Function

### Definition (Linear Model)

$$\hat{y}_i = \beta_0 + \beta_1 x_{i,1} + \beta_2 x_{i,2} + \cdots + \beta_m x_{i,m} = \beta_0 + \mathbf{x}_i^\top \boldsymbol{\beta}$$

- We need a way to assess how good our estimate  $\hat{y}_i$  approximates the expected output  $y_i$  given the current estimates of the coefficients  $\hat{\beta}_0, \dots, \hat{\beta}_m$ .
- Hence, we define a **loss function**  $\mathcal{L}(y, \hat{y})$ .



# Linear Model

## Loss Function

### Definition (Linear Model)

$$\hat{y}_i = \beta_0 + \beta_1 x_{i,1} + \beta_2 x_{i,2} + \cdots + \beta_m x_{i,m} = \beta_0 + \mathbf{x}_i^\top \boldsymbol{\beta}$$

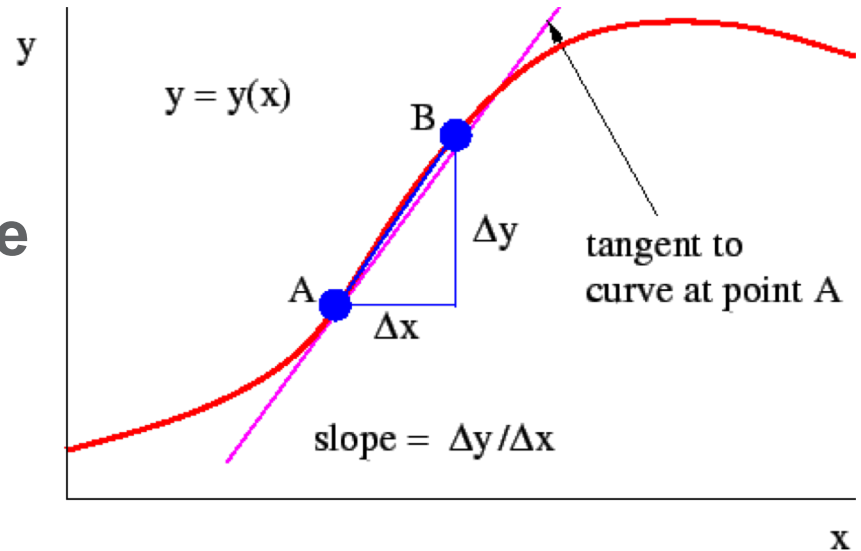
- If  $y_i$  is continuous such as "amount of body fat", a natural choice for the loss function is the **squared error**

$$\mathcal{L}(\mathbf{y}, \hat{\mathbf{y}}) = \frac{1}{2} \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \frac{1}{2} (\mathbf{y} - \hat{\mathbf{y}})^\top (\mathbf{y} - \hat{\mathbf{y}})$$

- $\mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})$  gives the total loss over the whole training set.
- We want to choose the coefficients  $\beta_0, \dots, \beta_m$  such that the total loss is **minimised**.
- Also referred to **residual sum of squares** in statistics.

# Derivative

- Consider a function  $y = f(x)$  with  $x, y \in \mathbb{R}$ .
- The derivative is denoted as  $f'(x)$  or  $\frac{d}{dx}f(x)$ .
- The **derivative** gives the **slope** of  $f(x)$  at the point  $x$ .
- It measures the sensitivity of  $f(x)$  to small changes in the input  $x$ .



# Partial Derivative and Gradient

- If the function has multiple inputs, we consider **partial derivatives**.
- The partial derivative  $\frac{\partial}{\partial x_i} f(\mathbf{x})$  measures how sensitive  $f(\mathbf{x})$  is to small changes in  $x_i$  alone at point  $\mathbf{x}$ .
- The **gradient** generalises the notion of derivative to the case where the derivative is with respect to a vector.
- The gradient of  $f(\mathbf{x})$  is the vector containing all the partial derivatives, denoted  $\nabla_{\mathbf{x}} f(\mathbf{x})$ .

# Function Minimisation

- We want to find parameters  $\beta_0, \boldsymbol{\beta}$  such that  $\mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})$  is minimized, i.e. solving

$$\arg \min_{\beta_0, \boldsymbol{\beta}} \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})$$
$$\arg \min_{\beta_0, \boldsymbol{\beta}} \frac{1}{2} \sum_{i=1}^n (y_i - \beta_0 - \mathbf{x}_i^\top \boldsymbol{\beta})^2$$

- $\mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})$  reaches its minimum at the point  $\beta_j$  if its **partial derivative** with respect to  $\beta_j$  is zero for all  $j = 0, \dots, m$ :

$$\frac{\partial}{\partial \beta_j} \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}}) = 0$$

# Ordinary Least Squares

## Estimation I

- Set the partial derivative to zero

$$\frac{\partial}{\partial \beta_j} \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_{i=1}^n x_{ij} (y_i - \beta_0 - \mathbf{x}_i^\top \boldsymbol{\beta}) = 0$$

- Gradient vector in matrix notation:

$$\nabla_{\boldsymbol{\beta}} \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}}) = -\mathbf{X}^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = \mathbf{0}$$

- Note:**  $\boldsymbol{\beta} = (\beta_0, \dots, \beta_m)^\top$  and the first column of  $\mathbf{X}$  contains only 1 to accommodate the intercept  $\beta_0$ , i.e.  $\mathbf{X}$  is a  $n \times m + 1$  matrix.



# Ordinary Least Squares

## Estimation II

### Definition (Ordinary Least Squares Estimate)

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

- The minimum of the loss function is unique.
- Estimates of the coefficients can be obtained in closed form and therefore no iterative optimisation is required.
- $\mathbf{X}$  must have full column rank, i.e.,  $\mathbf{X}^T \mathbf{X}$  is positive definite.
- Prediction is performed by

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_{i,1} + \hat{\beta}_2 x_{i,2} + \cdots + \hat{\beta}_m x_{i,m}$$

# Exercise I

- Open the notebook **01\_linear\_regression.ipynb**

# Logistic Regression

# Logistic Regression

- Consider a binary classification problem where  $y_i \in \{0, 1\}$ .
- If  $y_i = 1$ , the  $i$ -th sample belongs to the **positive class**, otherwise to the **negative class**.
- Create a model of the probability of samples  $x_i$  belonging to the positive class

$$\pi_i = P(y_i = 1 \mid x_{i,1}, \dots, x_{i,m})$$

- Remember that the linear model is define as

$$\eta_i = \beta_0 + \beta_1 x_{i,1} + \dots + \beta_m x_{i,m}$$

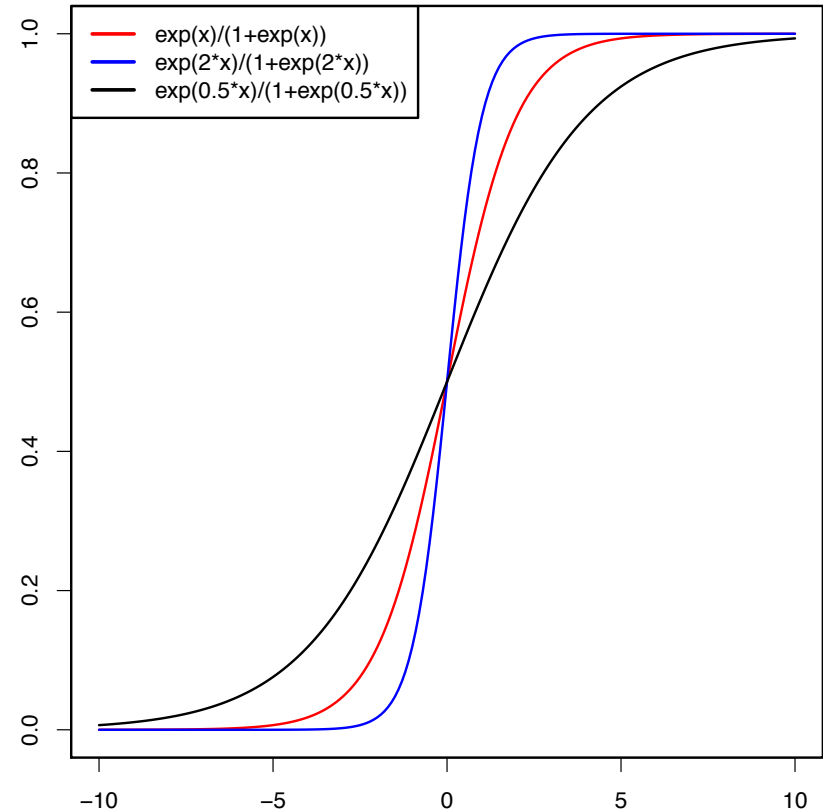
- **How to connect the probability  $\pi_i$  to the linear predictor  $\eta_i$ ?**

# Logistic Regression

## Response function

- The **logistic function**  $h(x)$  connects the probability  $\pi_i$  to the linear predictor  $\eta_i$

$$\pi_i = h(\eta_i) = \frac{\exp(\eta_i)}{1 + \exp(\eta_i)}$$



# Logistic Regression

## Loss function

### Definition (Likelihood function)

$$\mathcal{L}(\beta_o, \boldsymbol{\beta}; \mathbf{X}) = \prod_{i=1}^n P(y_i | \mathbf{x}_i) = \prod_{i=1}^n \pi_i^{y_i} (1 - \pi_i)^{1-y_i}$$

- We want **maximise** the probabilities across the whole training data.

### Definition (Maximum Likelihood Estimate; MLE)

$$\hat{\boldsymbol{\beta}} = \operatorname{argmax}_{\beta_o, \boldsymbol{\beta}} \log \mathcal{L}(\beta_o, \boldsymbol{\beta}; \mathbf{X})$$

# Logistic Regression

## Iterative Optimisation

- We proceed as before, by setting the gradient of  $\mathcal{L}(\beta_o, \boldsymbol{\beta}; \mathbf{X})$  to zero

$$\nabla_{\boldsymbol{\beta}} \mathcal{L}(\boldsymbol{\beta}) = \mathbf{X}^{\top} (\mathbf{y} - \boldsymbol{\pi}) = \mathbf{0}$$

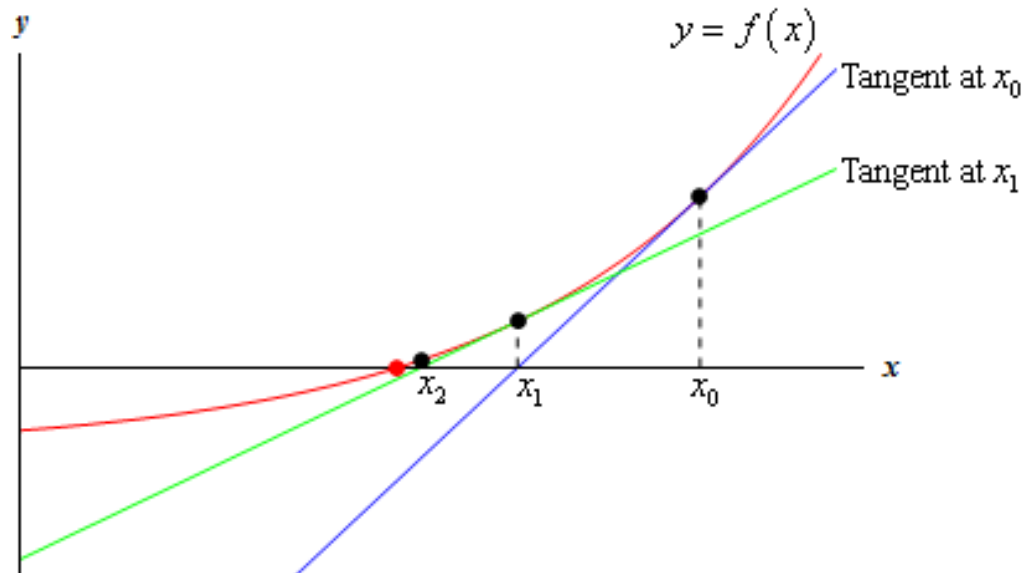
- However, we note that there is no closed form solution to finding parameters maximising the likelihood function.
- We need to **iteratively** search for the optimal set of parameters.

# Newton's method

48

- Given a function  $f(x)$  and its derivative  $f'(x)$ , **Newton's method** aims to find the value  $x$  that satisfies  $f(x) = 0$  by iteratively approximating the solution by calculating

$$x^{(t+1)} = x^{(t)} - \frac{f(x^{(t)})}{f'(x^{(t)})}$$





# Logistic Regression

## Iterative Optimisation

- We want to find  $\beta_o$  and  $\boldsymbol{\beta}$  for which

$$\nabla_{\boldsymbol{\beta}} \mathcal{L}(\boldsymbol{\beta}) = \mathbf{X}^T (\mathbf{y} - \boldsymbol{\pi}) = \mathbf{0}$$

- Using Newton's method, we find that the update step has the form

$$\boldsymbol{\beta}^{(t+1)} = \boldsymbol{\beta}^{(t)} - \left( \nabla_{\boldsymbol{\beta}}^2 \mathcal{L}(\boldsymbol{\beta}^{(t)}) \right)^{-1} \nabla_{\boldsymbol{\beta}} \mathcal{L}(\boldsymbol{\beta}^{(t)})$$

- Starting from an initial guess  $\boldsymbol{\beta}^{(0)}$ , we iteratively update our estimate of  $\boldsymbol{\beta}$  until convergence.
- Usually, the starting point is the zero vector,  $\boldsymbol{\beta}^{(0)} = \mathbf{0}$ .

# Second Derivative

## Hessian

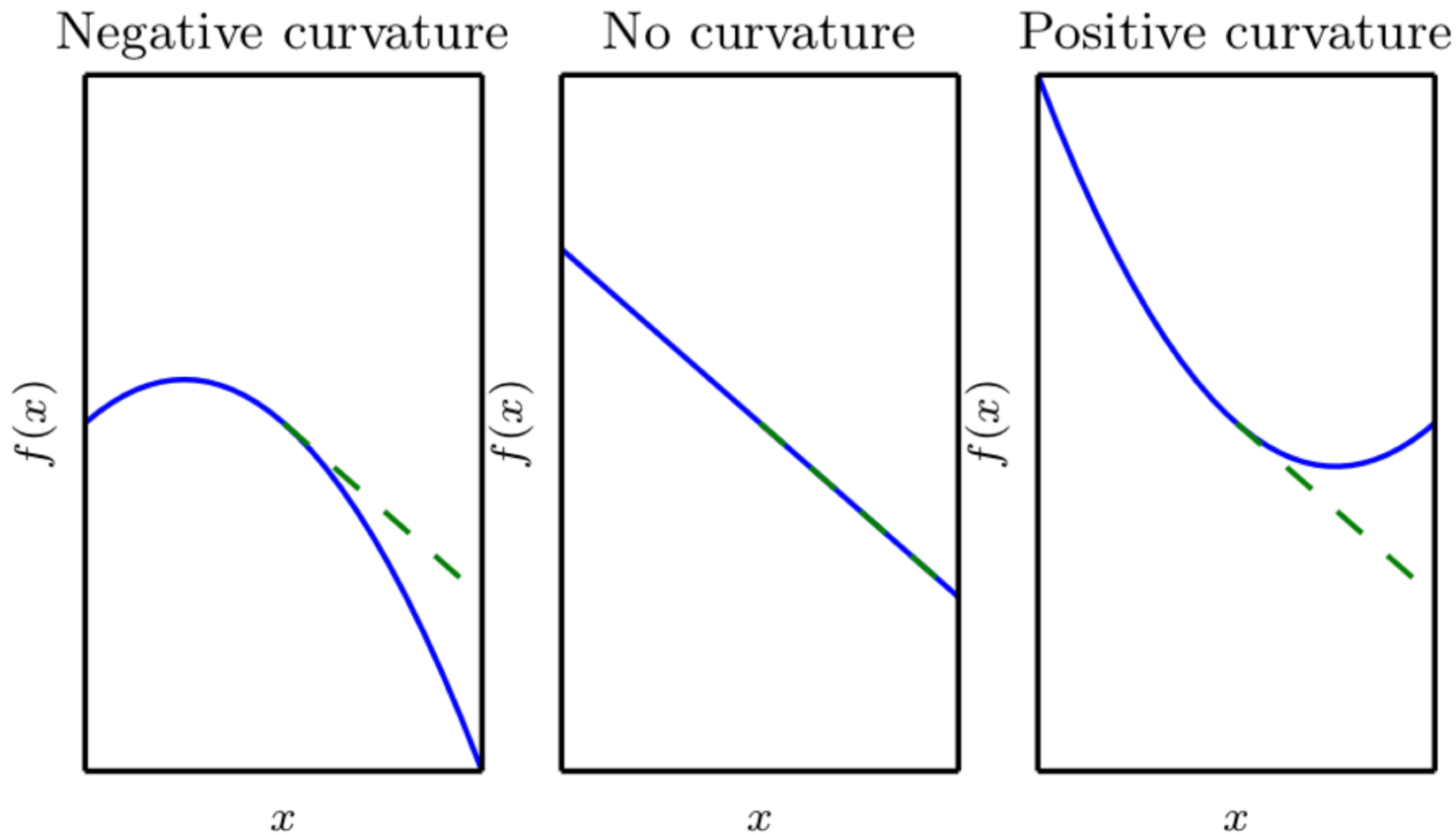
- The derivative of a derivative is called a **second-order derivative**.
- The **Hessian** matrix  $\mathbf{H} = \nabla_{\boldsymbol{\beta}}^2 \mathcal{L}(\boldsymbol{\beta}^{(t)})$  denotes a matrix of partial second-order derivatives

$$(\mathbf{H})_{i,j} = \frac{\partial^2}{\partial \beta_i \partial \beta_j} \mathcal{L}(\boldsymbol{\beta}^{(t)})$$

# Second Derivative

## Curvature

- The second-order derivative measures the **curvature** at a given point.

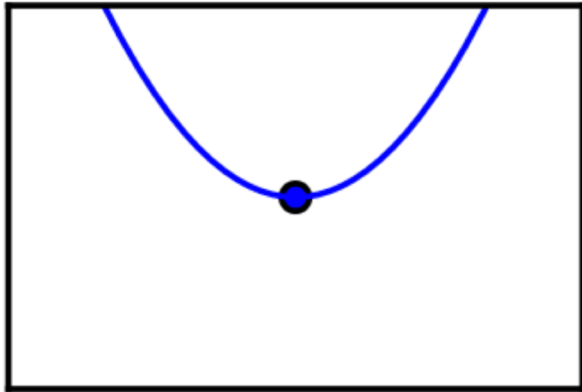


# Newton's Method

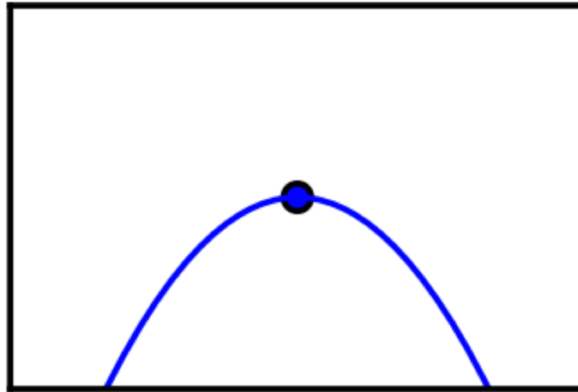
## Critical Points

52

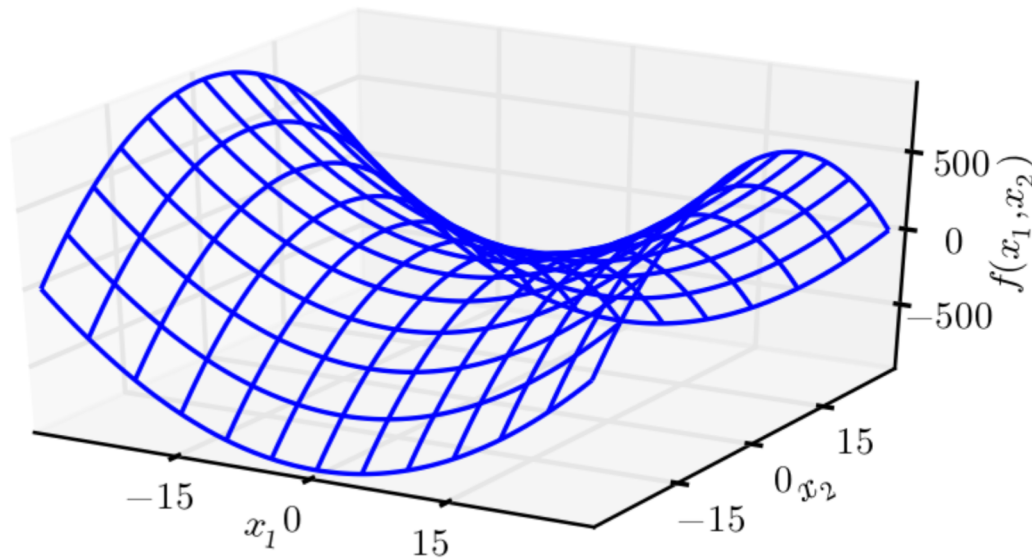
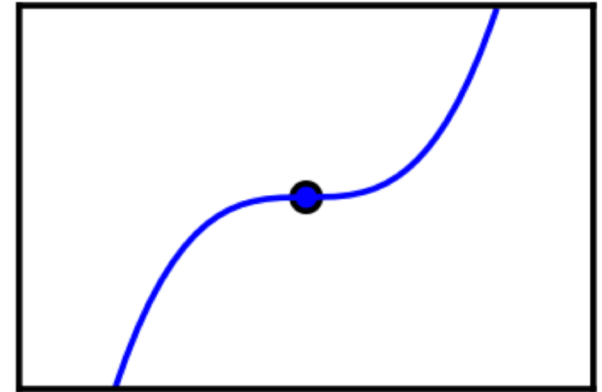
Minimum



Maximum



Saddle point



# Exercise II

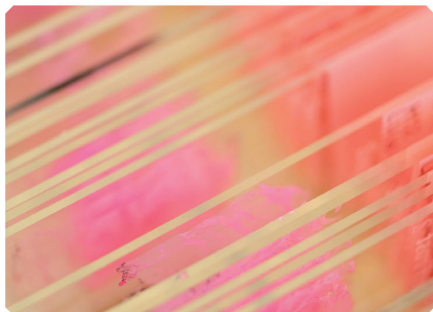
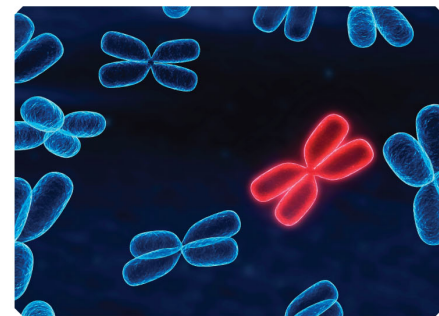
- Open the notebook **02\_logistic\_regression.ipynb**

ICR



Unrivalled  
track record

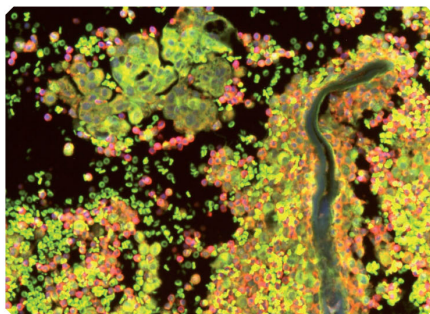
**ICR** The Institute of  
Cancer Research



Making the  
discoveries that  
defeat cancer



**ICR**



One of the world's  
most influential  
cancer research  
institutes