

# Introduction to Probabilistic Machine Learning

Linear Basis Function Models

Ralf Herbrich

# Overview

---

1. Linear Basis Function Models
2. Modelling Data
  - Modelling Text
  - Modelling Images
3. Linear Algebra
  - Vector Spaces
  - Linear Mappings and Matrices
  - Matrix Derivatives
4. Maximum A Posterior Learning and (Regularized) Least Squares

**Introduction to  
Probabilistic Machine  
Learning**

*Unit 6 – Linear Basis Function  
Models*

1. **Linear Basis Function Models**
2. Modelling Data
  - Modelling Text
  - Modelling Images
3. Linear Algebra
  - Vector Spaces
  - Linear Mappings and Matrices
  - Matrix Derivatives
4. Maximum A Posterior Learning and (Regularized) Least Squares

**Introduction to  
Probabilistic Machine  
Learning**

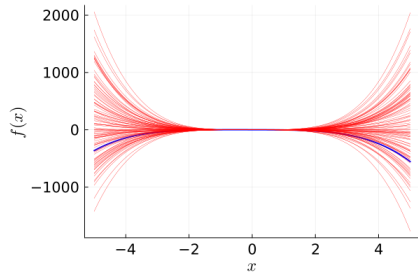
*Unit 6 – Linear Basis Function  
Models*

# Linear Basis Function Models

- **Linear Basis Function Models.** Given an input space  $\mathcal{X}$  and  $D$  basis functions  $\phi_j: \mathcal{X} \rightarrow \mathbb{R}$ , a linear basis function model is a function of the form
 
$$f(x; \mathbf{w}) = w_0 + w_1 \cdot \phi_1(x) + w_2 \cdot \phi_2(x) + \dots + w_D \cdot \phi_D(x)$$
- It's called a linear model, but the linearity is w.r.t.  $\mathbf{w}$  not  $x \in \mathcal{X}$ !
- **Examples:** 4 basis function ( $D = 4$ ) and 100 random parameters  $\mathbf{w}$ .

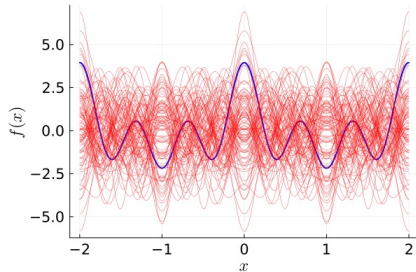
## Polynomial Basis

$$\phi_j(x) = x^j$$



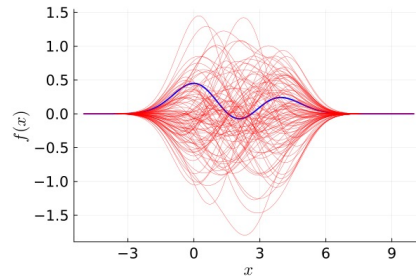
## Fourier Basis

$$\phi_j(x) = \cos(\pi j \cdot x)$$



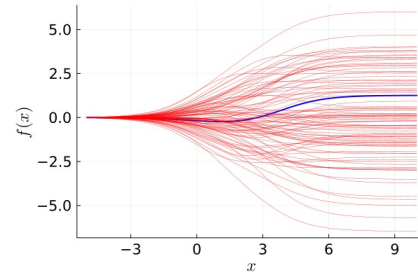
## Gaussian Basis

$$\phi_j(x) = \mathcal{N}(x; j, 1)$$



## Sigmoid Basis

$$\phi_j(x) = \frac{\exp(x - j)}{1 + \exp(x - j)}$$



# Maximum Likelihood and Least Squares

- **Gaussian Likelihood Model.** Let's assume that our observations  $y$  are obtained from adding zero-mean normally distributed noise to  $f(x; \mathbf{w})$

$$p(y|x, \mathbf{w}) = \mathcal{N}(y; f(x; \mathbf{w}), \sigma^2)$$

- **Maximum Likelihood Problem.** Noting that  $\log(\cdot)$  is a strictly monotonic function, we see that

$$\mathbf{w}_{\text{ML}} := \operatorname{argmax}_{\mathbf{w}} \prod_i p(y_i|x_i, \mathbf{w}) = \operatorname{argmax}_{\mathbf{w}} \sum_i \log(\mathcal{N}(y_i|f(x_i; \mathbf{w}), \sigma^2))$$

- Product and sum have the same maximum, but the objective function is numerically **much** more stable as a sum of logarithms of densities (why?)

- **Least Squares.** *The minimizer of the objective function  $\sum_i (y_i - f(x_i; \mathbf{w}))^2$  is  $\mathbf{w}_{\text{ML}}$ .*

- **Proof.** Looking at the normal density for a given  $(x_i, y_i)$  we see that

$$\log \left( \frac{1}{\sqrt{2\pi}\sigma} \exp \left( -\frac{(y_i - f(x_i; \mathbf{w}))^2}{2\sigma^2} \right) \right) = \log \left( \frac{1}{\sqrt{2\pi}\sigma} \right) - \frac{1}{2\sigma^2} \cdot \underbrace{(y_i - f(x_i; \mathbf{w}))^2}_{\text{only part dependent on } \mathbf{w}}$$

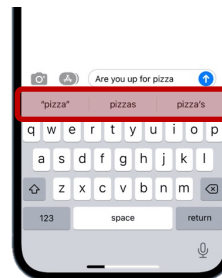
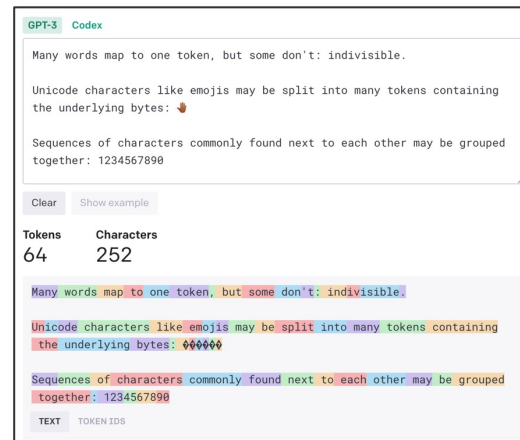
1. Linear Basis Function Models
2. **Modelling Data**
  - **Modelling Text**
  - Modelling Images
3. Linear Algebra
  - Vector Spaces
  - Linear Mappings and Matrices
  - Matrix Derivatives
4. Maximum A Posterior Learning and (Regularized) Least Squares

**Introduction to  
Probabilistic Machine  
Learning**

*Unit 6 – Linear Basis Function  
Models*

# Modelling Text

- Text can be modelled at three levels of granularity:
  1. **Letters** ( $\approx 10^2$  different letters in most alphabets)
  2. **Tokens** ( $10^3$  to  $10^4$  different tokens in most alphabets)
  3. **Words** ( $10^5$  to  $10^6$  different words in most languages)
- Modelling level of granularity depends on the application
  1. **Letters:** Compression algorithms for textual data
  2. **Tokens:** Sequence prediction for question-answering
  3. **Words:** Auto-correct function in smart keyboards



## Introduction to Probabilistic Machine Learning

Unit 6 – Linear Basis Function Models

# One-Hot Encoding

- Text is a sequence of text elements  $s_1, s_2, s_3, \dots$
- We often want to know the importance of each text element  $s_i$  to the target
  - **Example.** "This product shipped fast but was disappointing" has negative sentiment

$s_1$   $s_2$   $s_3$   $s_4$   $s_5$   $s_6$   $s_7$

Text element that is likely the expression of bad sentiment

- **One-Hot Encoding.** Given a dictionary  $S$  and a text element  $s \in S$ , a one-hot encoding  $\phi_{\text{OHE}}(s)$  is an  $|S|$ -dimensional unit vector indexed by the elements of  $S$  that consists of 0s in all dimensions with the exception of a single 1 in the dimension indexed by  $s$ .

- **Example (ctd).** If we assume the indices are  $s_1, s_2, \dots, s_7$  then

$$\phi_{\text{OHE}}(\text{"fast"}) = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \begin{array}{l} \leftarrow \text{This} \\ \leftarrow \text{product} \\ \leftarrow \text{shipped} \\ \leftarrow \text{fast} \\ \leftarrow \text{but} \\ \leftarrow \text{was} \\ \leftarrow \text{disappointing} \end{array}$$

**Introduction to  
Probabilistic Machine  
Learning**

Unit 6 – Linear Basis Function  
Models

- A linear model captures the *effect* of presence of a text element!



# Efficient One-Hot Encoding

- In practice, the feature vector for one-hot encoded data is *never* explicitly computed because

$$\mathbf{w}^T \boldsymbol{\phi}_{\text{OHE}}(s) = \sum_{i=1}^{|S|} w_i \cdot \phi_i(s) = w_{\text{idx}(s)}$$

$\mathcal{O}(|S|)$  (pointing to the sum)

$\mathcal{O}(\log_2 |S|)$  (pointing to  $\text{idx}(s)$ )

- All we need is the inverse function mapping of an actual text element  $s$

$$\text{idx}(s) := i \iff \phi_{\text{OHE},i}(s) = 1$$

- Two common techniques for encoding a whole text

- Variable-length text**  $s_1 s_2 \dots$ : Sum of all one-hot encoded vectors (“bag of words”)

$$\boldsymbol{\phi}_{\text{BOW}}(s_1 s_2 \dots) = \sum_j \boldsymbol{\phi}_{\text{OHE}}(s_j)$$

- Fixed-length text**  $s_1 s_2 \dots s_n$ : A stacked  $n \cdot |S|$  dimensional vector

$$\boldsymbol{\phi}_{\text{FL}}(s_1 s_2 \dots s_n) = \begin{bmatrix} \boldsymbol{\phi}_{\text{OHE}}(s_1) \\ \vdots \\ \boldsymbol{\phi}_{\text{OHE}}(s_n) \end{bmatrix}$$

Bag of words **cannot** learn  
“positional” effect of text elements

Introduction to  
Probabilistic Machine  
Learning

Unit 6 – Linear Basis Function  
Models

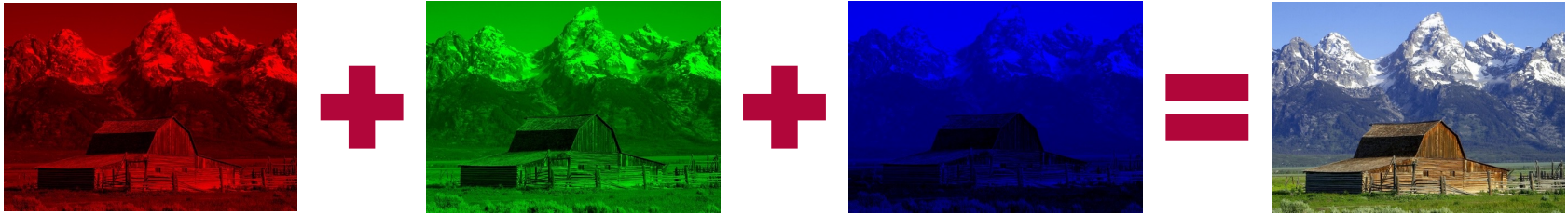
1. Linear Basis Function Models
2. **Modelling Data**
  - Modelling Text
  - **Modelling Images**
3. Linear Algebra
  - Vector Spaces
  - Linear Mappings and Matrices
  - Matrix Derivatives
4. Maximum A Posterior Learning and (Regularized) Least Squares

**Introduction to  
Probabilistic Machine  
Learning**

*Unit 6 – Linear Basis Function  
Models*

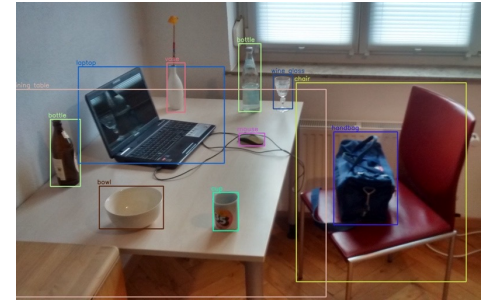
# Modelling Images

- **Raw image data.** *An image is a rectangular array of picture elements ("pixels") that consist of a triple of intensities of the base colors red, blue and green.*



- Images can also be modelled at three levels of granularity depending on application

1. **Pixels:** Image segmentation
2. (Non-overlapping) **patches:** Object recognition
3. **Whole image:** Image classification



# Image Content in Frequency and Location

- Image data contains signal (photon counts) at *fixed* locations in the image

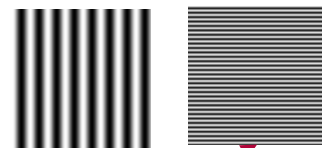
$$\phi_{\text{RGB}} \left( \text{Image} \right) = \begin{bmatrix} 127 \\ 0 \\ \vdots \\ 32 \\ 6 \end{bmatrix}$$

← Red value at location (1,1)  
 ← Green value at location (1,1)  
 ← Green value at location (800,600)  
 ← Blue value at location (800,600)

- Observation.** Recurring patterns are more visible in the frequency domain.
- Idea.** Discrete Fourier transform (DFT) to transform the image

$$\phi_{\text{DFT}}(x) = F_{\text{DFT}} \phi_{\text{RGB}}(x)$$

- In practice, there is a  $\mathcal{O}(\#\text{pixel} \cdot \log(\#\text{pixel}))$  algorithm for fast Fourier transform (FFT)

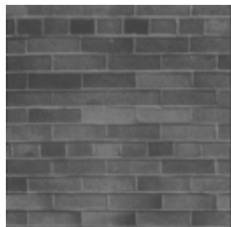


$$F_{\text{DFT}} = W_{\text{DFT}} \otimes W_{\text{DFT}}$$

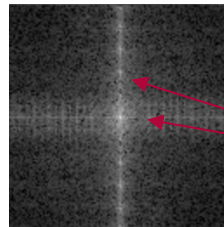
Kronecker product of 1D discrete Fourier transform  $W_{\text{DFT}}$

**Introduction to Probabilistic Machine Learning**

Unit 6 – Linear Basis Function Models



Raw image (in single-grey channel)

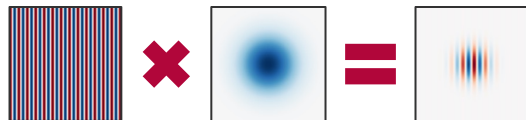


Fourier transform of image

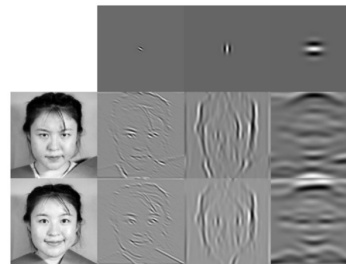
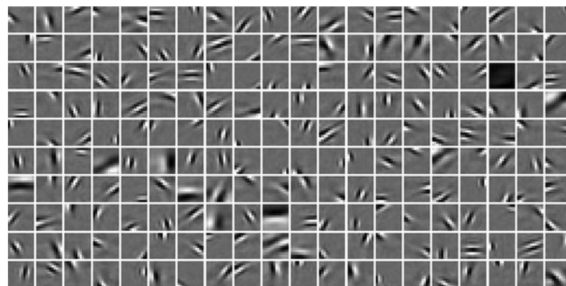
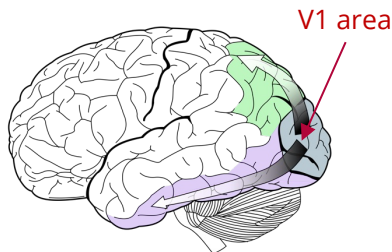
Horizontal and vertical patterns by single features

# Gabor Wavelets: Mixing Location and Frequency

- Both representations are extremes:
  - **Raw image:** Single features  $\phi_{\text{RGB},i}(x)$  at a location without frequency
  - **DFT image:** Single features  $\phi_{\text{DFT},i}(x)$  in frequency without a location
- **Gabor wavelets.** A *Gabor wavelet* is a set of basis functions that is obtained by multiplying a Fourier basis function with a Gaussian density.



- Combine both spatial and frequency information in image features
- Basis functions for sparse encoding by the brain in V1



Introduction to  
Probabilistic Machine  
Learning

Unit 6 – Linear Basis Function  
Models

Filters “implemented” in the V1  
(sparse coding!)

Bruno Olshausen & David Field (1997), [Sparse Coding with Overcomplete Basis Set: A Strategy Employed by V1?](#)



Dennis Gabor  
(1900 – 1978)

# Overview

---

1. Linear Basis Function Models
2. Modelling Data
  - Modelling Text
  - Modelling Images
3. **Linear Algebra**
  - **Vector Spaces**
  - Linear Mappings and Matrices
  - Matrix Derivatives
4. Maximum A Posterior Learning and (Regularized) Least Squares

**Introduction to  
Probabilistic Machine  
Learning**

*Unit 6 – Linear Basis Function  
Models*

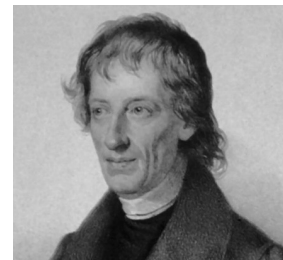
- **Vector Space.** A vector space is a set  $V$  that satisfies the following axioms: There exists a null element  $\mathbf{0} \in V$  such that for all  $\mathbf{u}, \mathbf{v}, \mathbf{w} \in V$  and scalars  $a, b \in \mathbb{R}$

1. **Associativity:**  $\mathbf{u} + (\mathbf{v} + \mathbf{w}) = (\mathbf{u} + \mathbf{v}) + \mathbf{w}$
2. **Commutativity:**  $\mathbf{u} + \mathbf{v} = \mathbf{v} + \mathbf{u}$
3. **Identity element (of vector addition):**  $\mathbf{v} + \mathbf{0} = \mathbf{v}$
4. **Inverse element:**  $\mathbf{v} + (-\mathbf{v}) = \mathbf{0}$
5. **Identity element (of scalar multiplication):**  $1 \cdot \mathbf{v} = \mathbf{v}$
6. **Distributivity (w.r.t. vector addition):**  $a \cdot (\mathbf{u} + \mathbf{v}) = a \cdot \mathbf{u} + a \cdot \mathbf{v}$
7. **Distributivity (w.r.t. scalar addition):**  $(a + b) \cdot \mathbf{v} = a \cdot \mathbf{v} + b \cdot \mathbf{v}$

- **Two Examples**

- **Points in space:**  $V = \mathbb{R}^d$  with  $\begin{pmatrix} u_1 \\ \vdots \\ u_d \end{pmatrix} + \begin{pmatrix} v_1 \\ \vdots \\ v_d \end{pmatrix} = \begin{pmatrix} u_1 + v_1 \\ \vdots \\ u_d + v_d \end{pmatrix}$  and  $a \cdot \begin{pmatrix} u_1 \\ \vdots \\ u_d \end{pmatrix} = \begin{pmatrix} au_1 \\ \vdots \\ au_d \end{pmatrix}$

- **Functions on  $\mathbb{R}$ :**  $V = \mathbb{R}^{\mathbb{R}}$  with  $f + g = x \mapsto f(x) + g(x)$  and  $a \cdot f = x \mapsto a \cdot f(x)$



Bernhard Bolzano  
(1781 – 1848)

**Introduction to  
Probabilistic Machine  
Learning**

Unit 6 – Linear Basis Function  
Models

# Linear Independence, Span and Basis

- **Linear combination.** Given a set  $v_1, v_2, \dots, v_n$  of a vector space  $V$ , a linear combination is defined by ( $a_1, a_2, \dots, a_n$  are called coefficients)

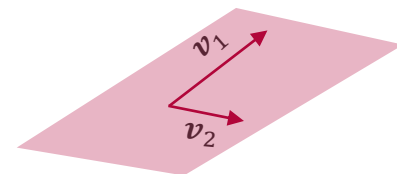
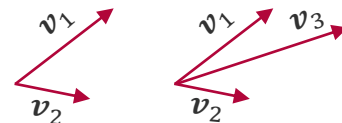
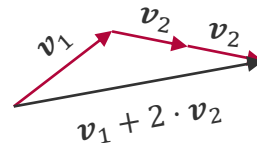
$$a_1 \cdot v_1 + a_2 \cdot v_2 + \dots + a_n \cdot v_n$$

- **Linear independence.** A set  $v_1, v_2, \dots, v_n$  is called linearly independent if no  $v_i$  can be written as a linear combination of  $v_1, \dots, v_{i-1}, v_{i+1}, \dots, v_n$ .
- **Span.** Given a set  $v_1, v_2, \dots, v_n$  of a vector space  $V$ , the span  $\text{span}(v_1, v_2, \dots, v_n)$  is the set of all linear combinations of  $v_1, v_2, \dots, v_n$ .

□ **Example:**

- $\text{span}(v_1)$  is the line through the origin along  $v_1$
- $\text{span}(v_1, v_2)$  is the plane through the origin along  $v_1$  and  $v_2$

- **Basis.** A subset  $b_1, b_2, \dots, b_n$  of a vector space  $V$  is called a basis if its span equals the whole vector space, that is,  $\text{span}(b_1, b_2, \dots, b_n) = V$ .
- **Dimensionality.** All bases of a vector space  $V$  have the same cardinality called the dimensionality of the vector space,  $\dim(V)$ .



**Introduction to  
Probabilistic Machine  
Learning**

Unit 6 – Linear Basis Function  
Models



# Scalar Products and Inner Product Spaces

- So far, a vector space is a set but there is no notion of **distance**!
  - In New York, the **walking distance** between two points is much longer than the **flying distance**!
- **Metric space.** A metric space is a vector space  $V$  together with a metric  $d: V \times V \rightarrow \mathbb{R}^+$  that has the following properties for all  $x, y, z$ :

$$d(x, y) = 0 \Leftrightarrow x = y$$

$$d(x, y) = d(y, x) \quad \leftarrow \text{symmetry}$$

$$d(x, y) \leq d(x, z) + d(z, y) \quad \leftarrow \text{triangle inequality}$$

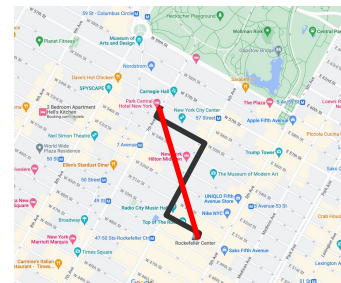
- **Inner product space.** An inner product space is a vector space together with an inner product function  $\langle \cdot, \cdot \rangle: V \times V \rightarrow \mathbb{R}^+$  that has the following properties for all  $x, y, z$ :

$$\langle x, x \rangle = 0 \Leftrightarrow x = \mathbf{0}$$

$$\langle x, y \rangle = \langle y, x \rangle \quad \leftarrow \text{symmetry}$$

$$\langle a \cdot x + b \cdot y, z \rangle = a \langle x, z \rangle + b \langle y, z \rangle \quad \leftarrow \text{linearity}$$

- Every inner product space is a metric space with  $d(x, y) = \|x - y\| = \sqrt{\langle x - y, x - y \rangle}$ !



Introduction to  
Probabilistic Machine  
Learning

Unit 6 – Linear Basis Function  
Models

1. Linear Basis Function Models
2. Modelling Data
  - Modelling Text
  - Modelling Images
3. **Linear Algebra**
  - Vector Spaces
  - **Linear Mappings and Matrices**
  - Matrix Derivatives
4. Maximum A Posterior Learning and (Regularized) Least Squares

**Introduction to  
Probabilistic Machine  
Learning**

*Unit 6 – Linear Basis Function  
Models*

# Linear Mappings and Matrices

- **Linear mapping.** A function  $f: V \rightarrow W$  is a linear mapping between vector space  $V$  and  $W$  if for any two vectors  $\mathbf{u}, \mathbf{v} \in V$  and any scalar  $c \in \mathbb{R}$ :

$$f(\mathbf{u} + \mathbf{v}) = f(\mathbf{u}) + f(\mathbf{v})$$

$$f(c \cdot \mathbf{u}) = c \cdot f(\mathbf{u})$$

- Addition and scalar multiplication can be applied before or after the map!
- It follows that  $f(\mathbf{0}) = f(0 \cdot \mathbf{v}) = 0 \cdot f(\mathbf{v}) = \mathbf{0}$ !

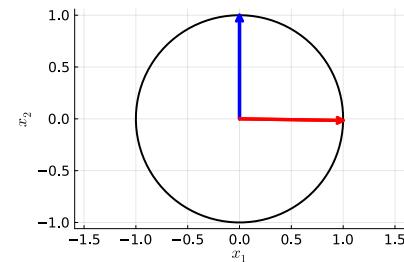
- **Theorem.** Every linear mapping  $f: V \rightarrow W$  between  $V$  and  $W$  of dimension  $n$  and  $m$  can be represented via a matrix multiplication with an  $m \times n$  matrix  $\mathbf{A}$ . The rank of  $\mathbf{A}$  is the dimensionality of the image of  $W$ , that is  $\text{rank}(\mathbf{A}) = \dim(W)$ .

- **Proof.** If  $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  is a basis for  $V$  and  $\{\mathbf{w}_1, \dots, \mathbf{w}_m\}$  is a basis for  $W$  then we know

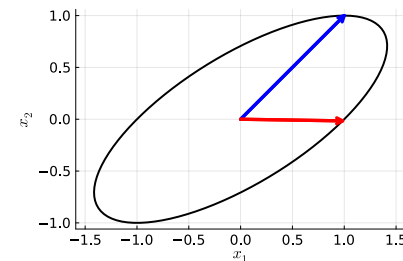
$$f(\mathbf{v}_j) = a_{1j} \cdot \mathbf{w}_1 + \dots + a_{mj} \cdot \mathbf{w}_m$$

$$f(\mathbf{v}) = f(c_1 \mathbf{v}_1 + \dots + c_n \mathbf{v}_n) = c_1 \cdot f(\mathbf{v}_1) + \dots + c_n \cdot f(\mathbf{v}_n) = \sum_{i=1}^m [\mathbf{A}c]_j \cdot \mathbf{w}_j$$

- For  $V = \mathbb{R}^n$  and  $W = \mathbb{R}^m$ , the columns of  $\mathbf{A}$  are the images of the basis vectors in  $\mathbb{R}^n$ .
- For any matrix  $\mathbf{A} \in \mathbb{R}^{n \times m}$ ,  $\text{rank}(\mathbf{A}) \leq \min(n, m)$ .



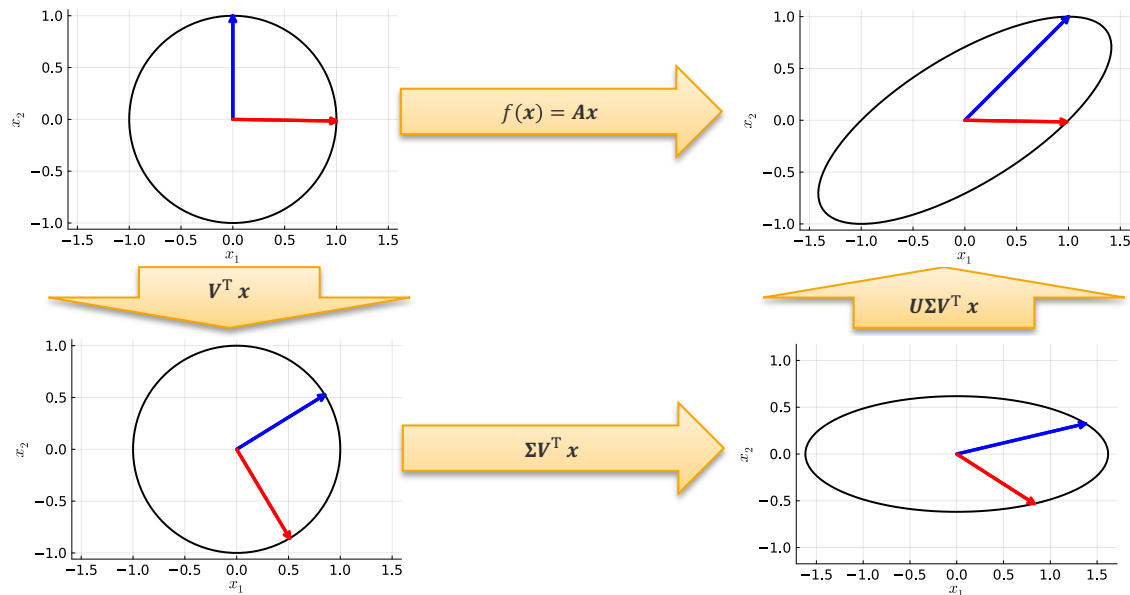
$$f(\mathbf{x}) = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix} \mathbf{x}$$



# Singular Value Decomposition

- **Singular Value Decomposition.** Any matrix  $A \in \mathbb{R}^{n \times m}$  has a decomposition into three matrices  $U \in \mathbb{R}^{n \times n}$ ,  $\Sigma \in \mathbb{R}^{n \times m}$  and  $V \in \mathbb{R}^{m \times m}$  such that  $UU^T$  and  $VV^T$  are the identity and  $\Sigma$  is only non-zero on diagonal elements

$$A = U\Sigma V^T$$



**Introduction to  
Probabilistic Machine  
Learning**

Unit 6 – Linear Basis Function  
Models

# Matrix Types and Properties

- **Square Matrix.** A matrix  $A \in \mathbb{R}^{n \times m}$  where  $m = n$  is called a square matrix. It parameterizes mappings of a space onto itself.
- **Symmetric Matrix.** A square matrix  $A$  is called symmetric if and only if  $A = A^T$ .
  - $A = U\Sigma V^T = (U\Sigma V^T)^T = V\Sigma U^T$  implies that  $V = U$  if all singular values are non-zero.
- **Diagonal Matrix.** A square matrix  $A$  is called diagonal if  $A_{ij} = 0$  for all  $i \neq j$ .
  - Geometrically, a diagonal matrix is an axis scaling (mapping).
  - A special diagonal matrix is  $A_{ii} = 1$  which is also called identity matrix.
- **Orthogonal Matrix.** A square matrix  $A$  is called orthogonal if  $AA^T = I$ .
  - Geometrically, an orthogonal matrix is a rotation and mirroring (mapping).
- **Positive semi-definite matrix.** A symmetric matrix  $A$  is called positive definite if and only if

$$\forall x \neq 0: x^T A x > 0$$

- If  $x = Uy$  then  $x^T A x = y^T U^T U \Sigma U U^T y = \sum_i \sigma_i y_i^2 > 0$ . Hence, positive definiteness means that all singular values are strictly positive, that is, no axis is inverted or removed.

**Introduction to  
Probabilistic Machine  
Learning**

Unit 6 – Linear Basis Function  
Models

# Determinant of a Matrix

- **Determinant.** The determinant,  $|A|$ , of a square matrix  $A \in \mathbb{R}^{n \times n}$  is the unique function that has the following three properties for any  $c \in \mathbb{R}$ :

$$|I| = 1$$

$$|[a_1, \dots, a_i, \dots, a_j, \dots, a_n]| = -|[a_1, \dots, a_j, \dots, a_i, \dots, a_n]|$$

$$|[a_1, \dots, c \cdot u + v, \dots, a_n]| = c \cdot |[a_1, \dots, u, \dots, a_n]| + |[a_1, \dots, v, \dots, a_n]|$$

- **Laplace Expansion.** If  $A_{[i,j]}$  is obtained from  $A$  by removing row  $i$  and column  $j$

$$\forall i: |A| = \sum_{j=1}^n (-1)^{i+j} A_{i,j} \cdot |A_{[i,j]}|$$

- Not useful for computation as  $O(n!)$  but useful for proofs!
- **Properties:** The determinant has the following two properties

$$|AB| = |A| \cdot |B|$$

$$|A^T| = |A|$$

Note that this implies for orthogonal matrices  $U$  that

$$1 = |I| = |UU^T| = |U| \cdot |U^T| = |U|^2 \Leftrightarrow |U| = \pm 1$$

- **Theorem:** The determinant  $|A|$  of a matrix  $A$  is the product of all singular values.

- **Proof.**  $|A| = |U\Sigma V^T| = |U| \cdot |\Sigma| \cdot |V^T| = |\Sigma| = \prod_i \sigma_i$
- **Corollary.** The determinant  $a$  is the volume of the parallelepiped of the new basis!



Pierre-Simon Laplace  
(1749 – 1827)

Unit 6 – Linear Basis Function Models

# Generalized Inverse of a Matrix



Eliakim Hastings Moore  
(1862 – 1932)



Sir Roger Penrose  
(1931 – )

Introduction to  
Probabilistic Machine  
Learning

Unit 6 – Linear Basis Function  
Models

- **Inverse.** The inverse  $A^{-1}$  of a full-rank square matrix  $A$  has the property that
$$A^{-1}A = AA^{-1} = I$$
- **Pseudo-Inverse (Moore-Penrose).** For any matrix  $A \in \mathbb{R}^{n \times m}$ , the pseudo-inverse  $A^+ \in \mathbb{R}^{m \times n}$  has the following properties (it is identical to  $A^{-1}$  when  $A$  is a full-rank and square matrix)

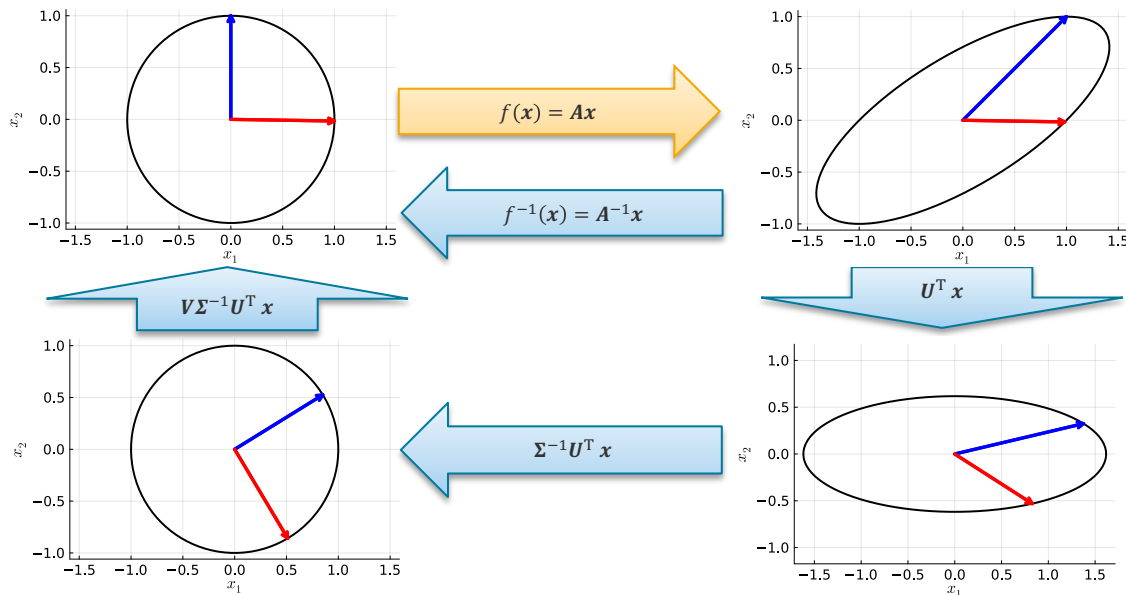
$$\begin{aligned}AA^+A &= A \text{ and } A^+AA^+ = A^+ \\(AA^+)^T &= AA^+ \text{ and } (A^+A)^T = A^+A\end{aligned}$$

- **Pseudo-Inverse Computation.** If a matrix either has full column or row rank
  - If  $(A^T A)^{-1}$  exists, then  $A^+ = (A^T A)^{-1} A^T$
  - If  $(A A^T)^{-1}$  exists, then  $A^+ = A^T (A A^T)^{-1}$
- **Least Squares Property.** For any matrix  $A$  and vector  $b$ , the minimizer of the least-squares difference between  $Ax$  and  $b$  is given by  $A^+b$ . In other words, for all  $x$ 
$$\|Ax - b\|^2 \geq \|AA^+b - b\|^2$$

# Generalized Inverse of a Matrix (SVD)

- But what if neither row or column rank are full?
  - Then we can use the singular value decomposition  $A = U\Sigma V^T$  because

$$A^+ = V\Sigma^{-1}U^T$$



**Introduction to  
Probabilistic Machine  
Learning**

*Unit 6 – Linear Basis Function  
Models*



# Maximum Likelihood Revisited

- **Maximum Likelihood Problem (in matrix notation).** If  $\Phi_{ij} = \phi_j(x_i)$

$$\mathbf{w}_{\text{ML}} := \operatorname{argmin}_{\mathbf{w}} \|\Phi \mathbf{w} - \mathbf{y}\|^2$$

- If we have more data points  $x_1, \dots, x_N$  than basis functions  $\phi_1, \dots, \phi_D$ , then  $\Phi^T \Phi$  has full rank and the minimizer can be computed from the Moore-Penrose inverse

$$\mathbf{w}_{\text{ML}} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y}$$

- **Computational complexity:**  $O(ND^2)$

1. The matrix multiplication  $\Phi^T \Phi$  is an  $O(ND^2)$  operation
2. The inverse  $(\Phi^T \Phi)^{-1}$  is an  $O(D^3)$  operation
3. The matrix product  $\Phi^T \mathbf{y}$  is an  $O(ND)$  operation
4. The matrix product  $\left[(\Phi^T \Phi)^{-1}\right] \cdot [\Phi^T \mathbf{y}]$  is an  $O(D^2)$  operation

- **Projection.** The matrix  $\mathbf{P} := \Phi(\Phi^T \Phi)^{-1} \Phi^T$  is a projection mapping of a target vector  $\mathbf{y}$  to its least-square approximation in the span of the basis functions!

- All projections have the property that  $\mathbf{P}\mathbf{P} = \mathbf{P}$ ! (proof)

**Introduction to  
Probabilistic Machine  
Learning**

Unit 6 – Linear Basis Function  
Models

# Cholesky Decomposition

- **Numerical Challenge:** Given a symmetric, positive-definite matrix  $A \in \mathbb{R}^{n \times n}$  and a vector  $b \in \mathbb{R}^n$  find the solution  $x$  such that

$$Ax = b$$

- **Naïve Solution:** Invert the matrix  $A$  and compute

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

- **Challenges:** If  $A$  has some singular values close to zero, this is numerically unstable!

- **Cholesky Decomposition:** Every positive-definitive matrix  $A \in \mathbb{R}^{n \times n}$  has a unique decomposition into a lower-triangular matrix  $L \in \mathbb{R}^{n \times n}$ ,  $L_{ij} = 0$  if  $j > i$

$$A = LL^T$$

- **Advantage:** Finding  $y$  such that  $Ly = b$  can be done in  $O(n^2)$  without an inverse simply using back-substitution (written as  $y = L \backslash b$ )!

- **Cholesky Solution:** Find the solution  $x$  for  $Ax = b$  is a two-step algorithm

1. Compute  $y = L \backslash b$
2. Compute  $x = L^T \backslash y$

$$Ax = b \Leftrightarrow LL^T x = b$$

$$\downarrow$$

$$Ly = b$$



André-Louis Cholesky  
(1875 – 1918)

$$\begin{pmatrix} L_{11} & 0 & 0 \\ L_{21} & L_{22} & 0 \\ L_{31} & L_{32} & L_{33} \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \\ y_3 \end{pmatrix} = \begin{pmatrix} b_1 \\ b \\ b_3 \end{pmatrix}$$

$$y_1 = \frac{b_1}{L_{11}}$$

$$y_2 = \frac{b_2 - L_{21} \cdot y_1}{L_{22}}$$

$$y_3 = \frac{b_3 - L_{31} \cdot y_1 - L_{32} \cdot y_2}{L_{33}}$$

# Overview

---

1. Linear Basis Function Models
2. Modelling Data
  - Modelling Text
  - Modelling Images
3. Linear Algebra
  - Vector Spaces
  - Linear Mappings and Matrices
  - **Matrix Derivatives**
4. Maximum A Posterior Learning and (Regularized) Least Squares

**Introduction to  
Probabilistic Machine  
Learning**

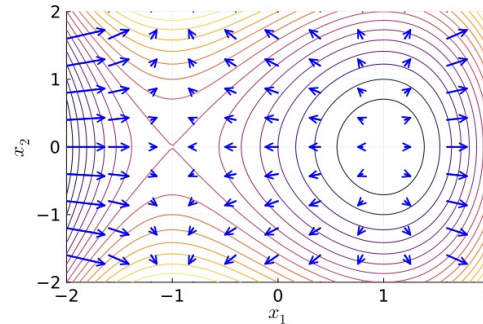
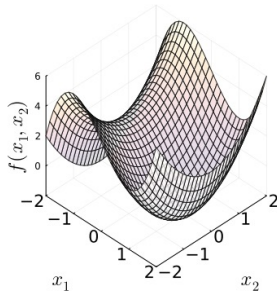
*Unit 6 – Linear Basis Function  
Models*

# Derivates of Vector-Valued Functions

- **Derivate of vector-valued functions.** Let  $g: \mathbb{R}^m \rightarrow \mathbb{R}^n$  be a vector-valued function. Then the  $m \times n$  matrix  $\frac{\partial g(x)}{\partial x}$  of derivatives is defined by

$$\frac{\partial g(x)}{\partial x} := \left[ \frac{\partial g_j(x)}{\partial x_i} \right]_{i,j=1}^{m,n} = \begin{bmatrix} \frac{dg_1(x)}{dx_1} & \dots & \frac{dg_n(x)}{dx_1} \\ \vdots & \ddots & \vdots \\ \frac{dg_1(x)}{dx_m} & \dots & \frac{dg_n(x)}{dx_m} \end{bmatrix}$$

- **Example.** Let  $g(x_1, x_2) = x_1^3 - 3x_1 + x_2^2$



**Introduction to  
Probabilistic Machine  
Learning**

Unit 6 – Linear Basis Function  
Models

# Derivates of Linear Mappings

- **Theorem.** Let  $f(x)$  be a linear mapping,  $f(x) = Ax + b$ . Then

$$\frac{\partial f(x)}{\partial x} = A^T$$

- **Theorem.** Let  $f(x)$  be a quadratic form,  $f(x) = x^T Ax$ . Then

$$\frac{\partial f(x)}{\partial x} = 2A^T x$$

- **Example.** Let us consider  $f(x) = \|Ax - b\|^2 = (Ax - b)^T (Ax - b)$

$$\begin{aligned} f(x) &= x^T A^T Ax - 2b^T Ax + b^T b \\ \left. \frac{\partial f(x)}{\partial x} \right|_{x=x^*} &= 2A^T Ax^* - 2A^T b = 0 \\ (A^T A)x^* &= A^T b \\ x^* &= \underbrace{(A^T A)^{-1} A^T}_{\text{Moore-Penrose Inverse}} b \end{aligned}$$

**Introduction to  
Probabilistic Machine  
Learning**

Unit 6 – Linear Basis Function  
Models

# Overview

---

1. Linear Basis Function Models
2. Modelling Data
  - Modelling Text
  - Modelling Images
3. Linear Algebra
  - Vector Spaces
  - Linear Mappings and Matrices
  - Matrix Derivatives
4. **Maximum A Posterior Learning and (Regularized) Least Squares**

**Introduction to  
Probabilistic Machine  
Learning**

*Unit 6 – Linear Basis Function  
Models*

# Maximum A Posteriori Learning

- **Gaussian Likelihood Model.** Let's assume that our observations  $y$  are obtained from adding zero-mean normally distributed noise

$$p(y|x, \mathbf{w}) = \mathcal{N}(y; f(x; \mathbf{w}), \sigma^2)$$

- **Gaussian Prior Model.** Let's assume a componentwise Gaussian over  $\mathbf{w}$

$$p(\mathbf{w}) = \prod_d \mathcal{N}(w_d; 0, \tau^2)$$

- **Maximum A-Posterior Problem.** Noting that  $\log(\cdot)$  is strictly monotonic

$$\mathbf{w}_{\text{MAP}} := \underset{\mathbf{w}}{\operatorname{argmax}} \underbrace{\prod_i p(y_i|x_i, \mathbf{w})}_{p(D|\mathbf{w})} \cdot \underbrace{\prod_d p(w_d)}_{p(\mathbf{w})} = \underset{\mathbf{w}}{\operatorname{argmin}} \frac{1}{2\sigma^2} \|\Phi \mathbf{w} - \mathbf{y}\|^2 + \frac{1}{2\tau^2} \mathbf{w}^T \mathbf{w}$$

$$\left. \frac{\partial \log(p(D|\mathbf{w}) \cdot p(\mathbf{w}))}{\partial \mathbf{w}} \right|_{\mathbf{w}=\mathbf{w}_{\text{MAP}}} = \frac{1}{\sigma^2} (\Phi^T \Phi \mathbf{w}_{\text{MAP}} - \Phi^T \mathbf{y}) + \frac{1}{\tau^2} \mathbf{w}_{\text{MAP}} = \mathbf{0}$$

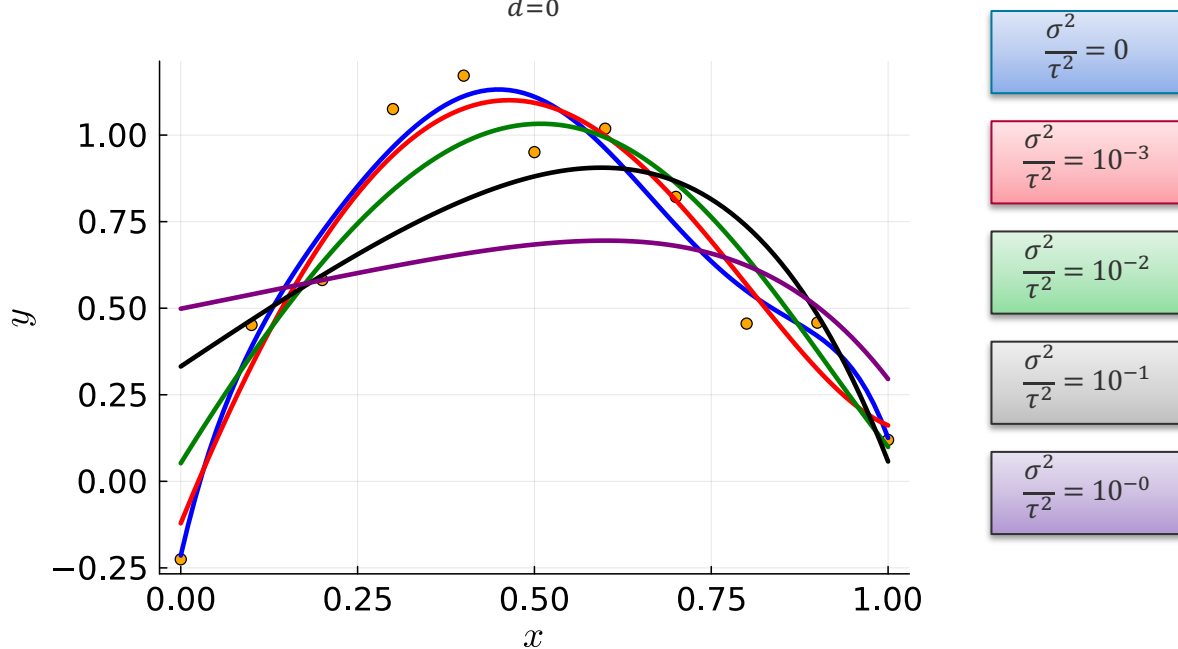
$$\mathbf{w}_{\text{MAP}} = \left( \Phi^T \Phi + \frac{\sigma^2}{\tau^2} \mathbf{I} \right)^{-1} \Phi^T \mathbf{y}$$

**Introduction to  
Probabilistic Machine  
Learning**

Unit 6 – Linear Basis Function  
Models

# Maximum A Posteriori Learning: Polynomial Regression

$$f(x; \mathbf{w}) = \sum_{d=0}^6 w_d \cdot x^d$$



**Introduction to  
Probabilistic Machine  
Learning**

*Unit 6 – Linear Basis Function  
Models*



## ■ Linear Basis Functions Models

- Non-linear prediction models can be formed with non-linear basis functions
- Linearity is in the parameters, *not* the input dimensions of the data

## ■ Modelling Data

- Both textual and image data can be represented at different levels of granularity
- Images should ideally be represented in terms of feature of location and frequency: Gabor wavelets/filters are excellent candidate functions for such features

## ■ Linear Mappings and Matrices

- All linear mappings can be expressed via matrix products
- The singular value decomposition is the rotation-scaling-rotation view on a mapping
- The inverse and determinant are critical for solving least-squares problems

## ■ Matrix Derivatives

- Matrix derivatives are a natural generalization of 1D-function derivatives
- Regularized least-square has closed-form solution using matrix derivatives

See you next week!