

# Applied Probabilistic Machine Learning

LINEAR, KERNEL, AND LOGISTIC REGRESSION

HUGUES RICHARD

*RichardH@rki.de*

IVAN TUNOV *Ivan.Tunov@student.hpi.uni-potsdam.de*

MF1 - GENOME COMPETENCE CENTER  
DACS

ROBERT KOCH INSTITUTE (RKI)  
HASSO PLATTNER INSTITUTE (HPI)

DECEMBER 19, 2024

(SLIDES COURTESY OF P. BENNER - BAM)

# REMINDER - SUPERVISED LEARNING

Training data  $\mathcal{D} = (\mathbf{x}_i, y_i)_{i=1, \dots, n}$

Find a function  $\hat{f}: \mathbf{x} \rightarrow \hat{f}(\mathbf{x}) = \mathbf{y}$

**Classification**

$\mathbf{y} \in \mathcal{C} = \{0, 1\}$  (binary)  
 $= \{1, \dots, C\}$  (multiclass)

Min. misclassification error:

$$\arg \min_{\hat{f}} \sum_{i=1}^n |\hat{f}(x_i) - y_i|$$

$$(\sum_{i=1}^n |\hat{f}(x_i) - y_i| = \sum_i \mathbb{I}_{\{\hat{f}(x_i) \neq y_i\}})$$

**Regression**

$\mathbf{y} \in \mathbb{R}$

Minimize a loss function  $\ell$

$$\sum_{i=1}^n \ell(\hat{f}(x_i) - y_i)$$

(e.g.  $\ell(x) = x^2, \ell(x) = |x|, \dots$ )

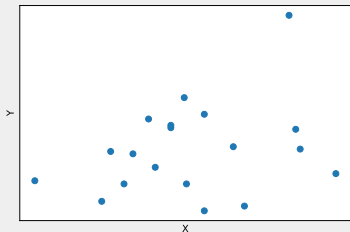
# LEARNING GOALS

- Understand linear regression and the probabilistic foundation between regression models
- Understand kernel regression when the relationship between features and outcome is not linear
- Understand logistic regression for classification.

# Linear Regression

# LINEAR REGRESSION

Let **Y** be the dependent variable (response variable) and **X** the independent variable (covariate, or predictor):



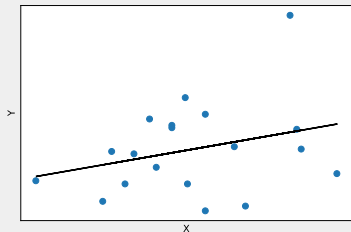
We assume the following model

$$\mathbf{Y} = f(\mathbf{X}) + \epsilon$$

where  $f$  is a linear function that models the expectation  $\mathbb{E}[Y | X]$ , and  $\epsilon$  is a noise term (e.g.  $\epsilon \sim \text{Normal}(0, \sigma^2)$ )

# LINEAR REGRESSION

Let **Y** be the dependent variable (response variable) and **X** the independent variable (covariate, or predictor):



We assume the following model

$$\mathbf{Y} = f(\mathbf{X}) + \epsilon$$

where  $f$  is a linear function that models the expectation  $\mathbb{E}[Y | X]$ , and  $\epsilon$  is a noise term (e.g.  $\epsilon \sim \text{Normal}(0, \sigma^2)$ )

# LINEAR REGRESSION

- We can also write  $\mathbf{Y} \sim \text{Normal}(f(\mathbf{X}), \sigma^2)$
- We assume no distribution for  $\mathbf{X}$
- We assume  $f$  is a linear function, i.e.

$$f(x) = ax + b$$

- How can we generate data  $(x_i, y_i)_i$  with this model?
  - ▶ For  $i = 1, \dots, n$ :
    - Select some value for  $x_i$
    - Draw  $\epsilon_i$  from  $\text{Normal}(0, \sigma^2)$
    - Compute  $y_i = f(x_i) + \epsilon_i$

# LINEAR REGRESSION - PARAMETER ESTIMATION

- In the Bayesian framework, parameters are estimated using the posterior distribution
- We want to know the probability of our hypothesis or parameters  $\theta = (a, b)$  given a set of  $n$  observations  $x = (x_i)_{i=1}^n$  and  $y = (y_i)_{i=1}^n$

- An estimate  $\hat{\theta}$  of our parameters  $\theta$  can be computed as the *maximum a posteriori (MAP) estimate*

$$\hat{\theta} = \arg \max_{\theta} \mathbb{P}(\theta \mid x, y)$$

- There are other choices, for instance the *posterior expectation*, which all have their justifications
- We use the MAP for linear regression, because it leads to a computationally simple solution



# LINEAR REGRESSION - PARAMETER ESTIMATION

- For a flat prior, the MAP is equivalent to the *maximum likelihood estimate (MLE)*, i.e.

$$\begin{aligned}\hat{\theta} &= \arg \max_{\theta} \mathbb{P}(\theta \mid x, y) \\ &= \arg \max_{\theta} \frac{\mathbb{P}(x, y \mid \theta) \mathbb{P}(\theta)}{\mathbb{P}(x, y)} \\ &= \arg \max_{\theta} \mathbb{P}(x, y \mid \theta) \mathbb{P}(\theta) \\ &= \arg \max_{\theta} \mathbb{P}(x, y \mid \theta)\end{aligned}$$

assuming  $\mathbb{P}(\theta)$  is constant<sup>1</sup>

- This result is not specific to linear regression models

---

<sup>1</sup>A uniform prior  $\mathbb{P}(\theta)$  is called *improper prior* when  $\theta$  is a continuous variable, because  $\mathbb{P}(\theta)$  does not integrate to one

- Furthermore, we have

$$\begin{aligned}\hat{\theta} &= \arg \max_{\theta} \mathbb{P}(x, y | \theta) \\ &= \arg \max_{\theta} \mathbb{P}(y | x, \theta) \mathbb{P}(x | \theta) \\ &= \arg \max_{\theta} \mathbb{P}(y | x, \theta)\end{aligned}$$

- In the last step we took advantage of the fact that the distribution of our covariates  $x$  does not depend on the parameters  $\theta$ , which are the slope and intercept of the linear function
- In fact, we do not have to assume a particular distribution for our covariates!

- Plugging in our normal distribution we arrive at

$$\begin{aligned}\hat{\theta} &= \arg \max_{\theta} \mathbb{P}(y_1 \dots y_n \mid x_1, \dots, x_n, \theta) \\&= \arg \max_{\theta} \prod_{i=1}^n \mathbb{P}(y_i \mid x_i, \theta) \\&= \arg \max_{\theta} \sum_{i=1}^n \log \mathbb{P}(y_i \mid x_i, \theta) \\&= \arg \max_{\theta} \sum_{i=1}^n \log \frac{1}{\sigma \sqrt{2\pi}} \exp \left\{ -\frac{(y_i - f(x_i))^2}{2\sigma^2} \right\} \\&= \arg \max_{\theta} \sum_{i=1}^n -(y_i - f(x_i))^2\end{aligned}$$

- The estimate

$$\begin{aligned}\hat{\theta} &= \arg \min_{\theta} \sum_{i=1}^n (y_i - f(x_i))^2 \\ &= \arg \min_{\theta} \sum_{i=1}^n (y_i - \hat{y}_i)^2\end{aligned}$$

is called the *ordinary least squares (OLS)* estimate

- It minimizes the squared error between our prediction  $\hat{y}_i$  and our observations  $y_i$
- In other words, it minimizes the squared residuals  $\epsilon_i = y_i - f(x_i)$

# LINEAR REGRESSION - GENERALIZATION

- For generalizing linear regression to multiple predictors, we first define

$$x = \begin{bmatrix} 1 \\ \tilde{x} \end{bmatrix}, \quad \theta = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix}$$

i.e.  $x$  is a vector where the first component is always 1

- This definition allows to write

$$\begin{aligned} f(x) &= b + a\tilde{x} \\ &= \theta_1 + \theta_2\tilde{x} \\ &= \begin{bmatrix} 1 \\ \tilde{x} \end{bmatrix}^\top \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix} \\ &= x^\top \theta \end{aligned}$$

# LINEAR REGRESSION - GENERALIZATION

- Adding additional predictors is now very simple

$$x = \begin{bmatrix} 1 \\ x^{(2)} \\ \vdots \\ x^{(p)} \end{bmatrix}, \quad \theta = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_p \end{bmatrix}$$

- The number of predictors / features is given by  $p$ , where the first predictor is  $(1, 1, \dots, 1)^\top$
- It follows that

$$\begin{aligned} f(x) &= x^\top \theta \\ &= \theta_1 + x^{(2)}\theta_2 + \dots + x^{(p)}\theta_p \end{aligned}$$

# LINEAR REGRESSION - NOTATION

- In general, we have  $n$  observations and  $p$  predictors
- For the  $i$ th observation  $(x_i, y_i)$ ,  $y_i$  is a scalar and  $x_i$  a vector

$$x_i = (1, x_i^{(2)}, \dots, x_i^{(p)})^\top$$

- We define the matrix

$$X = \begin{bmatrix} x_1^{(1)} & x_1^{(2)} & \dots & x_1^{(p)} \\ x_2^{(1)} & x_2^{(2)} & \dots & x_2^{(p)} \\ \vdots & \vdots & \ddots & \vdots \\ x_n^{(1)} & x_n^{(2)} & \dots & x_n^{(p)} \end{bmatrix} = \begin{bmatrix} 1 & x_1^{(2)} & \dots & x_1^{(p)} \\ 1 & x_2^{(2)} & \dots & x_2^{(p)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_n^{(2)} & \dots & x_n^{(p)} \end{bmatrix}$$

# LINEAR REGRESSION - NOTATION

- This notation allows us to write linear regression as

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} 1 & x_1^{(2)} & \dots & x_1^{(p)} \\ 1 & x_2^{(2)} & \dots & x_2^{(p)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_n^{(2)} & \dots & x_n^{(p)} \end{bmatrix} \begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_p \end{bmatrix} + \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix}$$

- Or in matrix notation simply as

$$y = X\theta + \epsilon$$

## Data matrix $X$

For a data matrix  $X \in \mathbb{R}^{n \times p}$ , rows will always correspond to observations and columns correspond to features. The first column is the vector  $(1, 1, \dots, 1)^\top$ . We always assume that  $X$  has full rank, i.e.  $\text{rank}(X) = \min(n, p)$



# LINEAR REGRESSION - OLS

If  $n > p$  and  $X^\top X$  has full rank we can use **ordinary least squared (OLS)** to estimate  $\theta$ :

$$\hat{\theta} = \arg \min_{\theta} \|\epsilon\|_2^2 = \arg \min_{\theta} \|y - X\theta\|_2^2$$

Differentiation with respect to  $\theta$  and solving for the roots leads to:

$$\begin{aligned} \Rightarrow \quad \hat{\theta} &= (X^\top X)^{-1} X^\top y \\ &= X^\top y \quad \text{if } X^\top X = I \end{aligned}$$

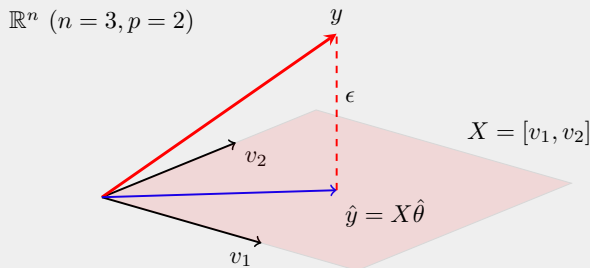
$X(X^\top X)^{-1}X^\top$  is called a projection matrix...

see exercise sheet for the derivation of the solution

# LINEAR REGRESSION - OLS PROJECTION

Let  $X\theta = v_1\theta_1 + v_2\theta_2 + \dots v_p\theta_p$ , where  $v_i$  denotes the  $i$ th column of  $X$

$$\hat{\theta} = \arg \min_{\theta} \|y - X\theta\|_2^2$$



$X(X^\top X)^{-1}X^\top y$  projects  $y$  onto the plane defined by the columns of  $X$

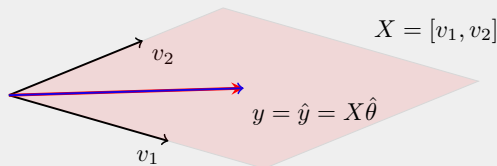
<sup>1</sup>[Hastie et al., 2009, Bishop, 2006]

# LINEAR REGRESSION - OLS PROJECTION

Let  $X\theta = v_1\theta_1 + v_2\theta_2 + \dots v_p\theta_p$ , where  $v_i$  denotes the  $i$ th column of  $X$

$$\hat{\theta} = \arg \min_{\theta} \|y - X\theta\|_2^2$$

$\mathbb{R}^n$  ( $n = 3, p = 2$ )



If  $y$  is already inside the plane, we obtain  $\epsilon = 0$

---

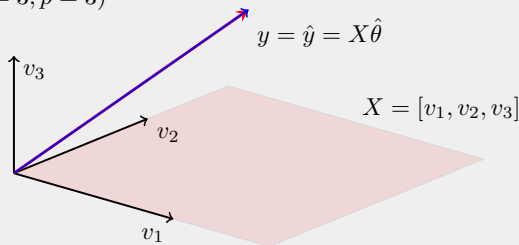
<sup>1</sup>[Hastie et al., 2009, Bishop, 2006]

# LINEAR REGRESSION - OLS PROJECTION

Let  $X\theta = v_1\theta_1 + v_2\theta_2 + \dots v_p\theta_p$ , where  $v_i$  denotes the  $i$ th column of  $X$

$$\hat{\theta} = \arg \min_{\theta} \|y - X\theta\|_2^2$$

$\mathbb{R}^n$  ( $n = 3, p = 3$ )



If  $p \geq n$  then  $\epsilon = 0$  and for  $p > n$  we have infinitely many solutions (assuming  $v_i$  are pairwise independent)

<sup>1</sup>[Hastie et al., 2009, Bishop, 2006]

- For  $p > n$  the OLS estimate

$$\hat{\theta} = \arg \min_{\theta} \|y - X\theta\|_2^2$$

has infinitely many solution  $\hat{\theta}$  such that  $\|y - X\hat{\theta}\|_2^2 = 0!$

- For  $p > n$  the OLS estimate

$$\hat{\theta} = \arg \min_{\theta} \|y - X\theta\|_2^2$$

has infinitely many solution  $\hat{\theta}$  such that  $\|y - X\hat{\theta}\|_2^2 = 0!$

- Which one should we choose?

- For  $p > n$  the OLS estimate

$$\hat{\theta} = \arg \min_{\theta} \|y - X\theta\|_2^2$$

has infinitely many solution  $\hat{\theta}$  such that  $\|y - X\hat{\theta}\|_2^2 = 0!$

- Which one should we choose?
- Remember our initial model

$$y = X\theta + \epsilon$$

and yet the estimate  $\hat{\theta}$  satisfies  $y = X\hat{\theta}$

# LINEAR REGRESSION - UNDERDETERMINED OLS

- For  $p > n$  the OLS estimate

$$\hat{\theta} = \arg \min_{\theta} \|y - X\theta\|_2^2$$

has infinitely many solution  $\hat{\theta}$  such that  $\|y - X\hat{\theta}\|_2^2 = 0!$

- Which one should we choose?
- Remember our initial model

$$y = X\theta + \epsilon$$

and yet the estimate  $\hat{\theta}$  satisfies  $y = X\hat{\theta}$

- Either  $\epsilon = 0$  or  $\hat{\theta}$  contains all the noise



# LINEAR REGRESSION - UNDERDETERMINED OLS

For instance, we could take that  $\theta$  with minimal length, i.e. the **minimum  $\ell_2$ -norm** solution<sup>2</sup>

$$\begin{aligned} \arg \min_{\theta} \quad & \|\theta\|_2^2 \\ \text{subject to} \quad & X\theta = y \end{aligned}$$

The solution is almost equivalent to the standard OLS solution, i.e.

$$\hat{\theta} = (X^\top X)^+ X^\top y$$

where  $(X^\top X)^+$  Moore-Penrose pseudoinverse<sup>3</sup> of  $X^\top X$ .

---

<sup>2</sup>**Common practice for training neural networks**

<sup>3</sup>The Moore-Penrose pseudoinverse of a matrix  $X$  is computed as follows: Let  $X = S\Sigma V^\top$  be the singular value decomposition of  $X$ , where  $\Sigma$  is a diagonal matrix containing the singular values.  $X^+ = S\Sigma^+ V^\top$  where  $\Sigma^+$  contains the reciprocal of all non-zero singular values.

## Ridge Regression

The ridge regression estimate is defined as

$$\hat{\theta}(\lambda) = \arg \min_{\theta} \|X\theta - y\|_2^2 + \lambda \|\theta\|_2^2$$

where  $\lambda$  is called the *regularization strength* or *penalty*. Note that  $\|\theta\|_2^2 = \sum_{i=2}^n \theta_i^2$ , i.e.  $\theta_1$  is not constrained

- There exists an analytical solution to the ridge estimate:

$$\hat{\theta}(\lambda) = (X^\top X + \lambda I)^{-1} X^\top y$$

- In the overparameterized case, for  $\lambda > 0$  we obtain  $\|\epsilon\|_2^2 > 0$

---

<sup>3</sup>Convex optimization: [Boyd and Vandenberghe, 2004]

# LINEAR REGRESSION - RIDGE REGRESSION

- For  $\lambda \rightarrow \infty$  the estimate  $\hat{\theta}(\lambda)$  converges to the componentwise regression estimator
- For  $\lambda \rightarrow 0$  the estimate  $\hat{\theta}(\lambda)$  converges to the minimum  $\ell_2$ -norm OLS solution<sup>4</sup>
- The penalty  $\lambda \|\theta\|_2^2$  can be interpreted as a Gaussian prior
- Ridge regression is useful when  $n < p$  and  $n \geq p$

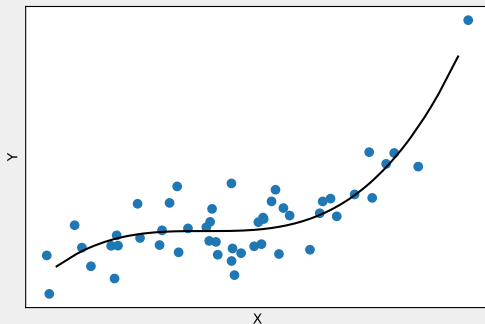
---

<sup>4</sup>  $A + \lambda I$  is invertible even for very small  $\lambda$ . In numerics,  $A + \lambda I$  is also used as a trick to ensure that a matrix is positive-definite.

# Kernel Regression

# POLYNOMIAL REGRESSION

- How can we change linear regression to model non-linear relations between  $X$  and  $Y$ ?



# REGRESSION IN FEATURE SPACE

Polynomial regression

$$\mathbf{Y} = \theta_1 + \theta_2 \mathbf{X} + \theta_3 \mathbf{X}^2 + \theta_4 \mathbf{X}^3 + \cdots + \epsilon,$$

More generally, we write

$$\mathbf{Y} = \phi(\mathbf{X})\theta + \epsilon,$$

where  $\phi : \mathbb{R}^p \rightarrow \mathbb{R}^{p'}$  is a **feature map** that maps points in  $p$ -dimensional input space into a  $p'$ -dimensional feature space, e.g.

$$\phi(\mathbf{X}) = (1, \mathbf{X}, \mathbf{X}^2, \mathbf{X}^3, \dots)$$

Basically linear (or ridge) regression in  $p'$ -dimensional feature space,  
but non-linear in input space

# KERNEL REGRESSION

- What if we do not know the exact set of features for our data?
- Can we simply test a large amount of possible features?
- Can we have more features than observations, i.e.  $n \leq p$ ?

Ridge regression in feature space:

$$\hat{\theta}(\lambda) = \arg \min_{\theta} \|\phi(X)\theta - y\|_2^2 + \lambda \|\theta\|_2^2$$

where  $\phi$  is applied to each row of  $X$ , i.e.  $\phi(X) \in \mathbb{R}^{n \times p'}$ .

Computationally expensive if  $p' \gg p$  and  $n \gg 1$ , assuming  $X$  is not sparse.

Reformulate the ridge regression estimate

$$\hat{\theta}(\lambda) = \arg \min_{\theta} \|\phi(X)\theta - y\|_2^2 + \lambda \|\theta\|_2^2$$

using **kernels**. Let  $\theta = \phi(X)^\top \eta$ , where  $\eta \in \mathbb{R}^n$  is a new parameter vector and  $\theta \in \text{span}(\phi(x_1), \dots, \phi(x_n)) \subset \mathbb{R}^p$ . It follows that

$$\begin{aligned}\hat{\eta}(\lambda) &= \arg \min_{\eta} \left\| \phi(X)\phi(X)^\top \eta - y \right\|_2^2 + \lambda \left\| \phi(X)^\top \eta \right\|_2^2 \\ &= \arg \min_{\eta} \|K\eta - y\|_2^2 + \lambda \eta^\top K\eta\end{aligned}$$

where  $K = \phi(X)\phi(X)^\top \in \mathbb{R}^{n \times n}$  is the **kernel matrix**.



## Definition: Kernel function

A function  $\kappa : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  is called a *kernel* if there exists a feature map  $\phi : \mathcal{X} \rightarrow \mathcal{F}$  such that

$$\kappa(x_i, x_j) = \phi(x_i)^\top \phi(x_j)$$

$K = (\kappa(x_i, x_j))_{x_i \in \mathcal{X}, x_j \in \mathcal{X}}$  is called the kernel matrix.

- $\mathcal{X}$  can be an arbitrary space, for instance DNA sequences
- $\kappa(x_i, x_j)$  is interpreted as a similarity measure in feature space
- Evaluating  $\kappa(x_i, x_j)$  does not always require to explicitly compute  $\phi(x)$
- Not having to map data into feature space is called the **kernel trick**

# EXAMPLE KERNELS

## ■ Linear kernel

$$\kappa(x_i, x_j) = x_i^\top x_j, \text{ where } \phi(x) = x$$

## ■ Polynomial kernel

$$\kappa(x_i, x_j) = (x_i^\top x_j + 1)^d$$

where  $d > 0$  is the degree. For  $\mathcal{X} = \mathbb{R}^2$  and  $d = 2$

$$\phi(x) = (1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, x_2^2, \sqrt{2}x_1x_2)^\top$$

## ■ Radial basis function (RBF) kernel

$$\kappa(x_i, x_j) = \exp\left(-\frac{\|x_i - x_j\|_2^2}{2\sigma^2}\right)$$

where the feature space has infinite dimensions

Let  $x_{\text{new}}$  denote the position where we would like to compute a prediction  $\hat{y}$

- Linear Regression

$$\hat{y} = \phi(x_{\text{new}})^\top \hat{\theta}$$

- Kernel Regression

$$\hat{y} = \sum_{i=1}^n \kappa(x_i, x_{\text{new}}) \hat{\eta}_i = \phi(x_{\text{new}})^\top \phi(X)^\top \hat{\eta}$$

which requires the full training set  $X = (x_i)_i \in \mathbb{R}^{n \times p}$ , where we simply used the definition  $\theta = \phi(X)^\top \eta$  to replace  $\hat{\theta}$  in the prediction of the linear regression model

# PARAMETERS AND HYPERPARAMETERS

- We call  $\theta$  and  $\eta$  the **parameters** of a (kernel) regression model
- The parameters of a kernel function (e.g.  $\sigma^2$  for the RBF kernel) or the regularization strength  $\lambda$  are also parameters of the model, but one step further up the hierarchy
- We call the parameters of a kernel function and the regularization strength **hyperparameters**
- In a Bayesian setting, the parameters control the likelihood function, whereas the hyperparameters parametrize the prior distribution

# KERNEL REGRESSION - PROS AND CONS

## Pros:

- Computationally efficient regression for high-dimensional feature spaces for moderate data sets
- Implicit regularization, i.e. only as many parameters as data points (but equivalent to minimum  $\ell_2$ -norm solution of standard regression)

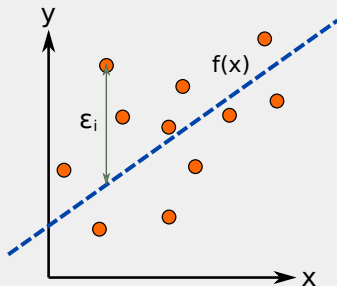
## Cons:

- Kernel matrix grows quadratically with number of samples
- $\theta \in \mathbb{R}^p \rightsquigarrow \eta \in \mathbb{R}^n$ , which creates dependencies between features
- Interpretation of parameters in feature space requires computation of  $\phi(X)^\top \eta$
- For infinite feature spaces  $\phi$  cannot be computed
- No feature selection possible ( $\ell_1$  penalty)

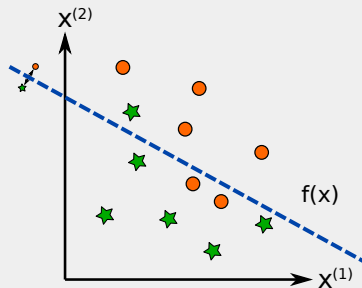
# Logistic Regression (Classification)

# LINEAR REGRESSION AND CLASSIFICATION

Linear regression



Logistic regression



$$y = X\theta + \epsilon, \quad y \in \mathbb{R} \quad y \stackrel{?}{=} \sigma(X\theta) + \epsilon, \quad y \in \{0, 1\}$$

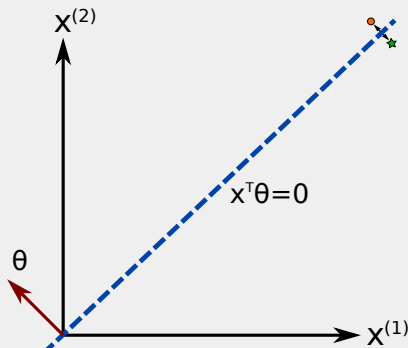
How is the hyperplane defined? What is  $\sigma$ ?

# DEFINING HYPERPLANES

- We use the properties of the dot product to define the separating hyperplane:

$$x^T \theta = \|x\| \|\theta\| \cos \angle$$

- For vectors  $x$  perpendicular to  $\theta$  we have  $\cos \angle = 0$

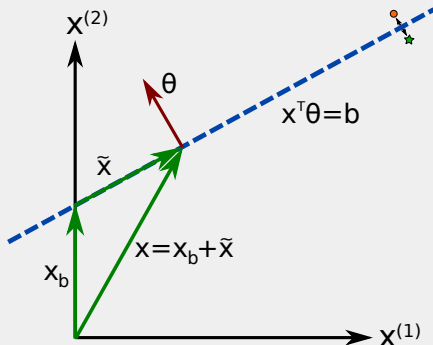




# DEFINING HYPERPLANES

- For hyperplanes with bias  $b$  we use  $x^\top \theta = b$

$$\begin{aligned}x^\top \theta &= (x_b + \tilde{x})^\top \theta \\&= \underbrace{x_b^\top \theta}_{=b} + \underbrace{\tilde{x}^\top \theta}_{=0}\end{aligned}$$



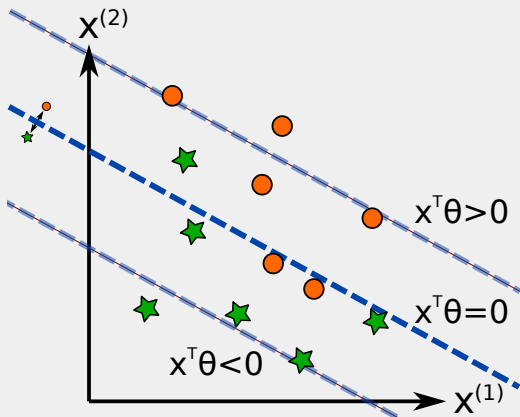
- Remember our convention:

$$x = \begin{bmatrix} 1 \\ x^{(2)} \\ \vdots \\ x^{(p)} \end{bmatrix}, \quad \theta = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_p \end{bmatrix}$$

- Hence, instead of  $x^\top \theta = b$  we can write  $x^\top \theta = 0$ , because  $\theta_1 = -b$

# SEPARATING HYPERPLANE

- $x^\top \theta > 0$  : predicting positive class
- $x^\top \theta < 0$  : predicting negative class



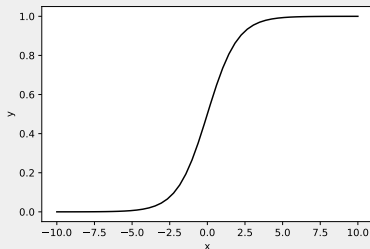
# LOGISTIC REGRESSION

- We convert  $x^\top \theta$  to probabilities

$$\mathbb{P}(Y = 1 \mid x) = \sigma(x^\top \theta)$$

- The function  $\sigma$  denotes the sigmoid function

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$



# LOGISTIC REGRESSION

- Given a training set  $(X, y)$  how do we estimate  $\theta$ ?
- Option 1: Minimizing squared error (similar to OLS)

$$\hat{\theta} = \arg \min_{\theta} \sum_{i=1}^n \left[ y_i - \sigma(x_i^{\top} \theta) \right]$$

Problem: Not convex!

- Remember how we justified OLS for linear models?
- Option 2: Maximum likelihood

$$\hat{\theta} = \arg \max_{\theta} \mathbb{P}(y | X, \theta)$$

# LOGISTIC REGRESSION

- What is the probability of  $(X, y)$ ?
- Remember a Bernoulli experiment (coin flip) with outcomes H (head) and T (tail)
- H is observed with probability  $p$
- T is observed with probability  $1 - p$
- The sequence HHTHT has probability

$$\mathbb{P}(\text{HHTHT}) = pp(1 - p)p(1 - p)$$

- Remember the following rule of thumb:

$\times$  = "and"

$+$  = "or"

# LOGISTIC REGRESSION

- For logistic regression, assume  $y = (1, 1, 0, 1)$ , hence

$$\mathbb{P}(1, 1, 0, 1 \mid X, \theta) = \sigma(x_1^\top \theta) \sigma(x_2^\top \theta) (1 - \sigma(x_3^\top \theta)) \sigma(x_4^\top \theta)$$

- Write it nicely in general form:

$$\mathbb{P}(y \mid X, \theta) = \prod_{i=1}^n \sigma(x_i^\top \theta)^{y_i} (1 - \sigma(x_i^\top \theta))^{1-y_i}$$

- Maximum likelihood

$$\begin{aligned} \hat{\theta} &= \arg \max_{\theta} \prod_{i=1}^n \sigma(x_i^\top \theta)^{y_i} (1 - \sigma(x_i^\top \theta))^{1-y_i} \\ &= \arg \max_{\theta} \sum_{i=1}^n y_i \log \sigma(x_i^\top \theta) + (1 - y_i) \log(1 - \sigma(x_i^\top \theta)) \end{aligned}$$



- Convex optimization problem, but must be solved numerically

# LEARNING GOALS

- Understand linear regression and the probabilistic foundation between regression models.
  - ▶ OLS is the Maximum a Posteriori / Maximum Likelihood of a linear relationship between input and target.
- Understand kernel regression when the relationship between features and outcome is not linear.
  - ▶ Kernel methods can account for non linear relationship with the kernel trick. It allows to understand many aspects of more complex models, such as neural networks
- Understand logistic regression for classification.
  - ▶ We can formulate binary classification as a regression problem on the separating hyperplane, with a probabilistic formulation. We will see many learning problem can be reformulated in a regression framework.



# REFERENCES

- 
- BISHOP, C. M. (2006).  
***Pattern Recognition and Machine Learning.***  
Springer.
- 
- BOYD, S. AND VANDENBERGHE, L. (2004).  
***Convex optimization.***  
Cambridge university press.
- 
- HASTIE, T., TIBSHIRANI, R., AND FRIEDMAN, J. (2009).  
***The elements of statistical learning: data mining, inference, and prediction.***  
Springer Science & Business Media.