

# Machine Learning

## MIRI Master

Lluís A. Belanche  
belanche@cs.upc.edu



Soft Computing Research Group  
*Departament de Ciències de la Computació* (Computer Science Department)  
Universitat Politècnica de Catalunya - Barcelona Tech

Spring Semester 2017-2018

**LECTURE 3: Linear regression and beyond**

# Linear regression and beyond

## Outline

1. Reminder of the regression framework
2. An example worked out
3. Leaping forward: basis functions and the SVD
4. Regularized least squares: ridge and the LASSO
5. Conclusions

# Linear regression and beyond

## Reminder of the regression framework

- The departing statistical **model** is

$$t_n = f(\mathbf{x}_n) + \varepsilon_n \quad \mathbf{x} \in \mathbb{R}^d, \quad t \in \mathbb{R}$$

where  $\varepsilon$  is a continuous r.v. such that  $\mathbb{E}[\varepsilon_n] = 0$  and  $\text{Var}[\varepsilon_n] = \sigma^2$

- Let's assume again that we further **model**  $\varepsilon_n \sim N(0, \sigma^2)$  and:

$$f(\mathbf{x}) \approx y(\mathbf{x}; \boldsymbol{\beta}) = \beta_0 + \sum_{i=1}^d \beta_i x_i = \boldsymbol{\beta}^\top \mathbf{x}$$

with  $\mathbf{x} = (1, x_1, \dots, x_d)^\top$  and  $\boldsymbol{\beta} = (\beta_0, \beta_1, \dots, \beta_d)^\top$

# Linear regression and beyond

## Reminder of the regression framework

Suppose we have an i.i.d. sample of  $N$  *labelled* observations  $\mathcal{D} = \{(\mathbf{x}_n, t_n)\}_{n=1, \dots, N}$ , where  $\mathbf{x}_n \in \mathbb{R}^d, t_n \in \mathbb{R}$

Therefore our **statistical model** is  $t_n \sim N(y(\mathbf{x}_n; \boldsymbol{\beta}), \sigma^2)$  or:

$$p(t_n | \mathbf{x}_n; \theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2} (t_n - \boldsymbol{\beta}^\top \mathbf{x}_n)^2\right),$$

with (unknown) **parameters**  $\theta := \{\beta_0, \beta_1, \dots, \beta_d, \sigma^2\}$ .

# Linear regression and beyond

## Reminder of the regression framework

Put  $\mathbf{t} = (t_1, \dots, t_N)^\top$  and  $X_{N \times (d+1)}$  the matrix of the  $\mathbf{x}_n$ . Define the **likelihood** as  $\mathcal{L}(\theta) := P(\mathbf{t}|\mathbf{X}; \theta)$

Let us maximize the “log-likelihood”:

$$\begin{aligned} l(\theta) &:= \ln \mathcal{L}(\theta) = \ln \prod_{n=1}^N p(t_n|\mathbf{x}_n; \theta) = \sum_{n=1}^N \ln p(t_n|\mathbf{x}_n; \theta) \\ &= -\frac{N}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{n=1}^N \left(t_n - \boldsymbol{\beta}^\top \mathbf{x}_n\right)^2 \\ &= -\frac{N}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} (\mathbf{t} - X\boldsymbol{\beta})^\top (\mathbf{t} - X\boldsymbol{\beta}) \\ &= -\frac{N}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \|\mathbf{t} - X\boldsymbol{\beta}\|^2 \end{aligned}$$

# Linear regression and beyond

## Reminder of the regression framework

$$\frac{\partial l}{\partial \beta} = -\frac{1}{2\sigma^2}(-2X^\top t + 2X^\top X\beta) = 0$$

$$\frac{\partial l}{\partial \sigma^2} = -\frac{N}{2\sigma^2} + \frac{1}{2\sigma^4}(t - X\beta)^\top (t - X\beta) = 0$$

Therefore,

$$\hat{\beta} = (X^\top X)^{-1} X^\top t$$

$$\hat{\sigma}^2 = \frac{1}{N}(t - X\hat{\beta})^\top (t - X\hat{\beta}) = \frac{1}{N}\|t - X\hat{\beta}\|^2$$

# Linear regression and beyond

## Reminder of the regression framework

- Note  $\hat{\sigma}^2 = R_{\text{emp}}(y_{\mathcal{D}})$ , which turns out to be a biased estimator of  $\sigma^2$ ; an unbiased estimator is:

$$\bar{\sigma}^2 = \frac{N}{N-d} \hat{\sigma}^2$$

- It is also known that  $\hat{\beta}$  is an unbiased estimator of  $\beta$  and

$$\text{Var}[\hat{\beta}] = (X^{\top} X)^{-1} \sigma^2$$

---

Which implies that  $\hat{\beta} \sim N(\beta, (X^{\top} X)^{-1} \sigma^2)$

# Linear regression and beyond

## Reminder of the regression framework

- The matrix  $X^\dagger := (X^\top X)^{-1}X^\top$  is known as the Moore-Penrose **pseudo-inverse** of  $X$
- It is the generalization of the notion of an inverse matrix to non-square matrices
- It has the property that  $X^\dagger X = I$  (although in general  $XX^\dagger \neq I$ ) (note however that both  $X^\dagger X$  and  $XX^\dagger$  are symmetric)



# Linear regression and beyond

## Reminder of the regression framework

**Theorem.** Let  $X_{N \times M}$  with  $N > M$ . If the column vectors of  $X$  are linearly independent, i.e., if  $\text{rank}(X) = M$ , then

1. the matrix  $X^\top X$  is symmetric and positive definite (p.d.) –in particular, it is non-singular
2. the least-squares problem

$$\min_{\beta \in \mathbb{R}^M} \|t - X\beta\|^2, \quad \text{has a unique solution}$$

3. this solution can be found solving the so-called Gauss' normal equations  $(X^\top X)\beta = X^\top t$  for  $\beta$

# Linear regression and beyond

## An example

The comet Tentax was discovered in 1968 and follows a quadratic orbit (elliptic, parabolic or hyperbolic) according to Kepler's laws.

The orbit has equation:

$$r = \frac{p}{1 - e \cos \varphi}$$

where  $p$  is a comet-specific coefficient,  $e$  is the comet's eccentricity (both unknown) and the  $(r, \varphi)$  pairs indicate measured positions (in polar coordinates centered at the Sun).

Astronomers have gathered a set of coordinates:

$\{(2,70, 48^\circ), (2,00, 67^\circ), (1,61, 83^\circ), (1,20, 108^\circ), (1,02, 126^\circ)\}$

**Goal:** estimate the two constants  $p, e$  based on measured data

# Linear regression and beyond

## An example

First we write the relation in linear form:  $r - (r \cos \varphi)e = p$  from which we arrive at the following system:

$$\left. \begin{array}{rcl} p + 1,806e & = & 2,70 \\ p + 0,782e & = & 2,00 \\ p + 0,196e & = & 1,61 \\ p - 0,371e & = & 1,20 \\ p - 0,600e & = & 1,02 \end{array} \right\}$$

which we express as:

$$X \cdot \beta = t \quad \text{or} \quad \begin{pmatrix} 1 & 1,806 \\ 1 & 0,782 \\ 1 & 0,196 \\ 1 & -0,371 \\ 1 & -0,600 \end{pmatrix} \cdot \begin{pmatrix} p \\ e \end{pmatrix} = \begin{pmatrix} 2,70 \\ 2,00 \\ 1,61 \\ 1,20 \\ 1,02 \end{pmatrix}$$

# Linear regression and beyond

## An example

Solving the normal equations  $(X^\top X)\beta = X^\top t$

as  $\hat{\beta} = (X^\top X)^{-1}X^\top t$

we obtain:

$$\hat{\beta} = \begin{pmatrix} p \\ e \end{pmatrix} \approx \begin{pmatrix} 1,454 \\ 0,694 \end{pmatrix}$$

# Linear regression and beyond

## Quality of the fit

- In statistics,  $-2l = -2 \ln \mathcal{L}$  is called the **deviance**
- In ML, this is the **square error**:

$$N \ln(2\pi\sigma^2) + \frac{1}{\sigma^2} \|t - X\hat{\beta}\|^2, \text{ that is to say } \|t - X\hat{\beta}\|^2$$

- A much better quantity to report is the **NRMSE**:

$$\text{NRMSE}(\hat{\beta}) = \sqrt{\frac{\|t - X\hat{\beta}\|^2}{(N - 1)\text{Var}[t]}}$$

---

In statistics,  $R^2 = 1 - \text{NRMSE}^2$  is the proportion of the (target) variability *explained* by the model

# Linear regression and beyond

## Leaping forward

We say that a model is **linear** if its parameters play a linear role in the model.

### Example

$$y(x; \boldsymbol{\beta}) = \beta_0 + \sum_{j=1}^d \beta_j x^j = \beta_0 + \beta_1 x + \dots + \beta_d x^d, \quad x \in \mathbb{R}$$

is a polynomial on  $x$  but a linear model!

# Linear regression and beyond

## Leaping forward

A simple but powerful idea is the introduction of **basis functions**:

$$y(\mathbf{x}; \mathbf{w}) = w_0 + \sum_{j=1}^M w_j \phi_j(\mathbf{x}) = \sum_{j=0}^M w_j \phi_j(\mathbf{x}) = \mathbf{w}^\top \boldsymbol{\phi}(\mathbf{x})$$

where  $\phi_0(\mathbf{x}) = 1$ .

$\boldsymbol{\phi}(\mathbf{x}) = (1, \phi_1(\mathbf{x}), \dots, \phi_M(\mathbf{x}))^\top$  and  $\mathbf{w} = (w_0, w_1, \dots, w_M)^\top$ .

---

The basis function expansion above is still a **linear model**.

# Linear regression and beyond

## Leaping forward

**Example:**

$$y(x; \mathbf{w}) = w_0 + \sum_{j=1}^M w_j \phi_j(x) = \sum_{j=0}^M w_j \phi_j(x) = \mathbf{w}^\top \boldsymbol{\phi}(x)$$

where  $\phi_0(x) = \sqrt{T^{-1}}$  and  $M$  is even

$$\phi_j(x) = \begin{cases} \sqrt{2T^{-1}} \sin(a_j x), & a_j = (j+1)\pi T^{-1} \quad \text{if } j \text{ is odd} \\ \sqrt{2T^{-1}} \cos(a_j x), & a_j = j\pi T^{-1} \quad \text{if } j \text{ is even} \end{cases}$$

This  $y$  is a truncated Fourier series in  $[0, T]$ . If  $T$  is known, it is a linear model for  $x$ ; what if  $T$  is unknown (then a parameter)?



# Linear regression and beyond

## Leaping forward

Define  $\mathbf{t} = (t_1, \dots, t_N)^\top$  the vector of targets

$\Phi_{N \times (M+1)}$  the matrix of the  $\phi(\mathbf{x}_n)$

where  $\Phi_{ij} = \phi_j(\mathbf{x}_i)$ ,  $i = 1, \dots, N, j = 0, \dots, M$ .

$$\Phi = \begin{pmatrix} 1 & \phi_1(\mathbf{x}_1) & \phi_2(\mathbf{x}_1) & \dots & \phi_M(\mathbf{x}_1) \\ 1 & \phi_1(\mathbf{x}_2) & \phi_2(\mathbf{x}_2) & \dots & \phi_M(\mathbf{x}_2) \\ \dots & \dots & \dots & \ddots & \vdots \\ 1 & \phi_1(\mathbf{x}_N) & \phi_2(\mathbf{x}_N) & \dots & \phi_M(\mathbf{x}_N) \end{pmatrix}$$

# Linear regression and beyond

## Leaping forward

Let us maximize the new log-likelihood:

1. The Gauss' normal equations are  $(\Phi^\top \Phi)w = \Phi^\top t$

2. The solution of which is

$$\hat{w} = (\Phi^\top \Phi)^{-1} \Phi^\top t = \Phi^\dagger t$$

and

$$\hat{\sigma}^2 = \frac{1}{N} (t - \Phi \hat{w})^\top (t - \Phi \hat{w}) = \frac{1}{N} \|t - \Phi \hat{w}\|^2$$

# Linear regression and beyond

## Singular Value Decomposition

The direct computation of the pseudo-inverse of  $\Phi$  has two major drawbacks:

1. When  $M$  is large,  $(\Phi^\top \Phi)$  is a large  $(M + 1) \times (M + 1)$  matrix; then the required inverse  $(\Phi^\top \Phi)^{-1}$  can be costly
2. If  $(\Phi^\top \Phi)$  is singular –or close to– then the required inverse  $(\Phi^\top \Phi)^{-1}$  can be impossible –or numerically delicate

# Linear regression and beyond

## Singular Value Decomposition

**Theorem.** Every matrix  $X_{N \times M}$  can be expressed as:

$$X = U \Delta V^\top$$

$U_{N \times N}, V_{M \times M}$  are **orthogonal** matrices ( $U^\top U = I_N, V^\top V = I_M$ )  
 $\Delta_{N \times M}$  is a **diagonal** matrix

---

The columns of  $U$  are the **eigenvectors** of  $XX^\top$   
The columns of  $V$  are the **eigenvectors** of  $X^\top X$

# Linear regression and beyond

## Singular Value Decomposition

1. Let  $\text{rank}(X) = r \leq \min(N, M)$ . Then exactly  $r$  elements  $\lambda_k$  in the diagonal of  $\Delta$  are strictly positive; the remaining elements are 0
2. These  $\lambda_k > 0$  are called the **singular values** and correspond to the square roots of the positive **eigenvalues** of  $XX^\top$  (same as  $X^\top X$ )
3. Sometimes an “economy size” decomposition is delivered:

If  $X$  is  $N \times M$  with  $N > M$ , then only the first  $M$  columns of  $U$  are given and  $\Delta$  is  $M \times M$

# Linear regression and beyond

## SVD example

$$X = \begin{pmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \\ 7 & 8 \end{pmatrix}$$

$$U = \begin{pmatrix} -0,1525 & -0,8226 & -0,3945 & -0,3800 \\ -0,3499 & -0,4214 & 0,2428 & 0,8007 \\ -0,5474 & -0,0201 & 0,6979 & -0,4614 \\ -0,7448 & 0,3812 & -0,5462 & 0,0407 \end{pmatrix}$$

$$\Delta = \begin{pmatrix} 14,2691 & 0 \\ 0 & 0,6268 \\ 0 & 0 \\ 0 & 0 \end{pmatrix}$$

$$V = \begin{pmatrix} -0,6414 & 0,7672 \\ -0,7672 & -0,6414 \end{pmatrix}$$

# Linear regression and beyond

## SVD economy size example

$$X = \begin{pmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \\ 7 & 8 \end{pmatrix}$$

$$U = \begin{pmatrix} -0,1525 & -0,8226 \\ -0,3499 & -0,4214 \\ -0,5474 & -0,0201 \\ -0,7448 & 0,3812 \end{pmatrix}$$

$$\Delta = \begin{pmatrix} 14,2691 & 0 \\ 0 & 0,6268 \end{pmatrix}$$

$$V = \begin{pmatrix} -0,6414 & 0,7672 \\ -0,7672 & -0,6414 \end{pmatrix}$$

# Linear regression and beyond

## The SVD for least squares

Given the least-squares problem

$$\min_{\mathbf{w} \in \mathbb{R}^M} \|\mathbf{t} - X\mathbf{w}\|^2$$

the solution can be obtained with the SVD as:

1. Compute the economy size SVD of  $X = U\Delta V^\top$
2. Solve for  $\mathbf{w}$  as  $\hat{\mathbf{w}} = V \text{diag}(\lambda_k^{-1}) U^\top \mathbf{t}$ , where only the  $\lambda_k > 0$  are considered



# Linear regression and beyond

## Regularized least squares

The maximum likelihood framework can yield unstable parameter estimates, especially when:

1. The explanatory variables are highly correlated
2. There is an insufficient number of observations ( $N$ ) relative to the number of predictors (basis functions  $M + 1$  or dimensions  $d + 1$ )

# Linear regression and beyond

## Regularized least squares

In the context of regression with Gaussian noise (square error), it is quite common to penalize the parameter vector:

1. Define the **penalized empirical error** as:

$$R_{\text{emp}}(y(\cdot; \mathbf{w})) := \|\mathbf{t} - \Phi \mathbf{w}\|^2 + \lambda \|\mathbf{w}\|^2, \quad \lambda > 0$$

2. Set the derivative w.r.t.  $\mathbf{w}$  to 0:  $(-2\Phi^\top \mathbf{t} + 2\Phi^\top \Phi \mathbf{w}) + 2\lambda \mathbf{w} = 0$
3. Therefore,  $\hat{\mathbf{w}} = (\Phi^\top \Phi + \lambda I)^{-1} \Phi^\top \mathbf{t}$

# Linear regression and beyond

## Regularized least squares

- This is known as Tikhonov or  $L_2$  **regularization** in ML
- Perhaps it is best known as **ridge regression** in statistics, where it is usually explained as a “penalized log-likelihood”
- It can also be derived from Bayesian statistics arguments
- Advantages:
  1. Pushing the length of the parameter vector  $\|\mathbf{w}\|$  to 0 allows the fit to be under explicit control with the regularization parameter  $\lambda$
  2. The matrix  $\Phi^\top \Phi$  is positive semi-definite (p.s.d.); therefore  $\Phi^\top \Phi + \lambda I$  is guaranteed to be p.d. (hence non-singular),  $\forall \lambda > 0$

# Linear regression and beyond

## Regularized least squares

Yes, nice, but ...

- How to do the **explicit control** on the fit?
  - regularization permits the specification of models that are more complex than needed because it limits the effective complexity
  - instead of trial-and-error on complexity, we can set a large complexity and adjust the  $\lambda$

# Linear regression and beyond

## Regularized least squares

Yes, **very** nice, **but** ...

- How to set the value of  $\lambda$ ? Using LOOCV, because
  - in this case  $\lambda$  is a very forgiving parameter (we usually perform a log search)
  - there is a closed efficient formula for the LOOCV (for **linear models**)

# Linear regression and beyond

## Regularized least squares

1. Choose a (large) set of values  $\Lambda$

2. For every  $\lambda \in \Lambda$ ,

a) Solve for  $\hat{\mathbf{w}} = (\Phi^\top \Phi + \lambda I)^{-1} \Phi^\top \mathbf{t}$

b) Compute the “hat matrix”  $H := \Phi \Phi^\dagger = \Phi (\Phi^\top \Phi + \lambda I)^{-1} \Phi^\top$

c) Compute the LOOCV of  $y(\cdot) = \hat{\mathbf{w}}^\top \phi(\cdot)$  in  $\mathcal{D}$  as

$$\text{LOOCV}(y) = \frac{1}{N} \sum_{n=1}^N \left( \frac{t_n - \hat{\mathbf{w}}^\top \phi(\mathbf{x}_n)}{1 - h_{nn}} \right)^2$$

3. Choose the model with the lowest LOOCV

# Linear regression and beyond

## Regularized least squares

A very popular method is GCV (**generalized cross-validation**):

$$\text{GCV}(y) = \frac{1}{N} \frac{\sum_{n=1}^N \left( t_n - \hat{\mathbf{w}}^\top \phi(\mathbf{x}_n) \right)^2}{\left( 1 - \frac{\text{Tr}(H)}{N} \right)^2}$$

which is a more stable computation for the LOOCV

Note that  $\lambda$  is needed to compute both  $\hat{\mathbf{w}}$  and  $H$

# Linear regression and beyond

## LASSO regression

The LASSO (Least Absolute Shrinkage and Selection Operator) regression is  $L_1$ -regularized Linear regression

The choice for the regularizer is  $\|\mathbf{w}\|_1$  and we get:

$$R_{\text{emp}}(y(\cdot; \mathbf{w})) = \|\mathbf{t} - \Phi \mathbf{w}\|^2 + \tau \|\mathbf{w}\|_1, \quad \tau > 0$$

which turns out to be equivalent to

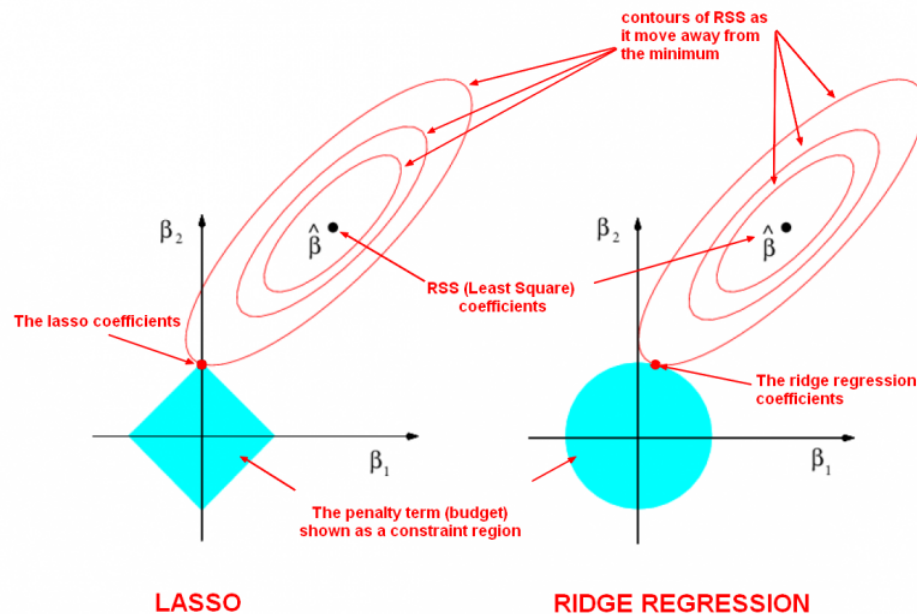
$$R_{\text{emp}}(y(\cdot; \mathbf{w})) = \|\mathbf{t} - \Phi \mathbf{w}\|^2, \text{ subject to } \|\mathbf{w}\|_1 \leq \tau$$



# Linear regression and beyond

## LASSO regression

In ridge regression, as the penalty  $\lambda$  is increased, all coefficients are reduced while still remaining non-zero; in the LASSO, increasing the  $\tau$  penalty causes more and more of the coefficients to be driven to zero



As  $d$  increases, the multidimensional diamond has an increasing number of corners, and so it is highly likely that some coefficients will be set equal to zero. Hence, the LASSO performs **shrinkage** and therefore **feature selection**

# Linear regression and beyond

## LASSO regression

- The LASSO loss function is no longer quadratic, but is still convex
- The LASSO is a special **quadratic programming** (QP) problem, for which the **Least Angle Regression** (LARS) procedure is used
- It exploits the special structure of the problem, and provides an efficient way to compute the solutions for all possible values of  $\tau > 0$  (**regularization path**)

# Linear regression and beyond

## Conclusions

We have introduced **linear models** as linear combinations of non-linear **basis functions** (BF)

### ADVANTAGES:

1. We can represent non-linear functions of the data using linear fitting techniques; we have the freedom to choose the form of the BFs
2. The fit can be under tight explicit control by regularization
3. The computations can be very efficient, no need to refit for LOOCV
4. Interpretability of the model is rather high

# Linear regression and beyond

## Conclusions

### LIMITATIONS:

The most important weak point is in the basis functions!

1. Many interesting basis functions scale very poorly with dimension (polynomials, Fourier series, splines, ...)
2. Our BFs are not flexible (they are independent of the data)
3. As a consequence, their number may be very high, which in turn leads to instability (because of low significance of the coefficients)

# Linear regression and beyond

## Conclusions

The solution is to develop **basis functions with parameters** that

1. ... scale well with dimension (inner products, distances, ...)
2. ... are data dependent (because of the parameters)
3. As a consequence, their number may be much lower (and the coefficients be significant)
4. Unfortunately, the new parameters play a non-linear role in the model: their optimization is plagued with local optima

# Machine Learning

## Syllabus

1. Introduction to Machine Learning
2. Theoretical issues (I): regression
3. Linear regression and beyond
4. Theoretical issues (II): classification
5. Generative classifiers
6. Discriminative classifiers

7. Clustering: k-means and E-M
8. Learning with kernels (I): The SVM
9. Learning with kernels (II): Kernel functions
10. Artificial neural networks (I): Delta rule, MLP-1
11. Artificial neural networks (II): MLP-2, RBF
12. Artificial neural networks (III): DL and CNNs
13. Ensemble methods: Random Forests
14. Advanced topics and frontiers