

## References

- This slides are based on material of prof. Marcello Restelli

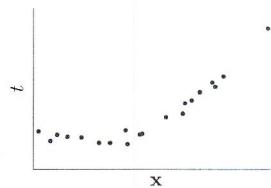
- *Pattern Recognition and Machine Learning*, Bishop
  - ▶ Chapter 1 (1.1, 1.2, 1.3)
  - ▶ Chapter 3 (3.1, 3.3)



## What is regression?

- Learn an **approximation** of function  $f(x)$  that maps input  $x$  to a continuous output  $t$  from a dataset  $\mathcal{D}$

$$\mathcal{D} = \{(x, t)\} \Rightarrow t = f(x)$$

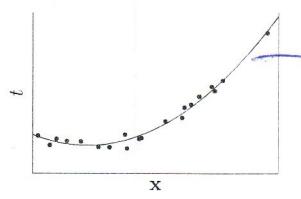


## What is regression?

- Learn an **approximation** of function  $f(x)$  that maps input  $x$  to a continuous output  $t$  from a dataset  $\mathcal{D}$

$$\mathcal{D} = \{(x, t)\} \Rightarrow t = f(x)$$

- ▶ How do we model  $f$ ?
- ▶ How do we evaluate our approximation?
- ▶ How do we optimize our approximation?



the function that generates the points (plus some gaussian noise)

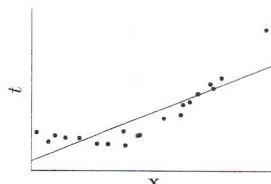
Note: we'll always have some noise, even in processes that are not stochastic



## Linear Regression

- In linear regression,  $f(x)$  is modeled with linear functions
  - ▶ Linear models can be easily explained
  - ▶ A linear regression problem can be solved **analytically**
  - ▶ Linear functions can be extended to model also non-linear relationships
  - ▶ More sophisticated methods are based on linear regression

We can solve the optimization problem analytically. With this property we're sure that if we didn't achieve good results it's not because of the optimization algorithm (this is very big)



## Linear Regression: model

- The simplest linear model can be defined as:

$$y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{j=1}^{D-1} w_j x_j = \mathbf{w}^T \mathbf{x}$$

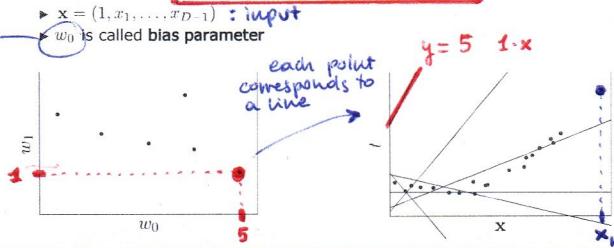
needed to model functions that have an offset  $w_0$  on the input

Suppose  $D=2$  (dimension). The dataset becomes:

$$\mathbf{D} = \begin{bmatrix} 1, x_1 \\ 1, x_2 \\ \vdots \\ 1, x_n \end{bmatrix}$$

In this case the hypothesis space is the space of all the possible values of the two weights:  $w_0, w_1$ . One point in this hypothesis space (i.e. one couple  $(w_0, w_1)$ ) represents a model:  $y = w_0 + w_1 x$ . How can we evaluate all of the models (points) and select the best one?

Remember:  
dimension =  $(1 + \# \text{ variables})$   
for the constant term ( $w_0$ )



Machine Learning

Daniele Lolacano

## Linear Regression: loss function and optimization

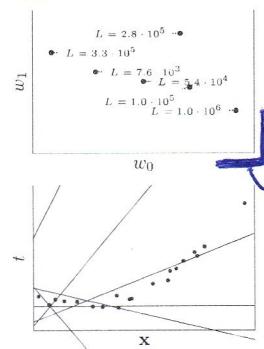
- A convenient error loss function is the sum of squared errors (SSE):

$$L(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (y(x_n, \mathbf{w}) - t_n)^2$$

► the sum in  $L$  is also called residual sum of squares (RSS) and can be written as the sum of residual errors:

$$RSS(\mathbf{w}) = \|\mathbf{\epsilon}\|_2^2 = \sum_{i=1}^N \epsilon_i^2$$

► closed-form optimization of  $L$  can be easily obtained



Given a point  $(w_0, w_1)$  we can compute  $L(w)$  by computing the sum of the squared errors that our approximation is doing w.r.t. the actual target value. The  $\frac{1}{2}$  term is just for convenience (when we'll optimize).

To each point  $(w_0, w_1)$  corresponds a value  $L(w)$ . The best model is the one for which  $L(w)$  is smaller.

Machine Learning

Daniele Lolacano

Linear Models and Basis Functions :  
Improve linear models with basis functions

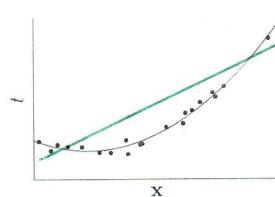
Machine Learning

Daniele Lolacano

- A linear combination of the input variables is not enough to model data...
- ... but we just need a regression model that is linear in the parameters
- We can define a model using non-linear basis functions:

$$y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{j=1}^{M-1} w_j \phi_j(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x})$$

$$\phi(\mathbf{x}) = (1, \phi_1(\mathbf{x}), \dots, \phi_{M-1}(\mathbf{x}))^T$$



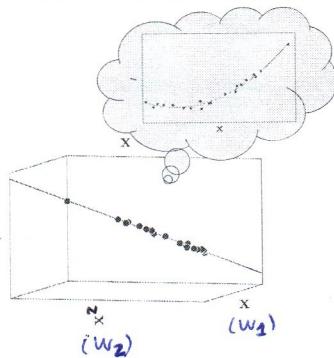
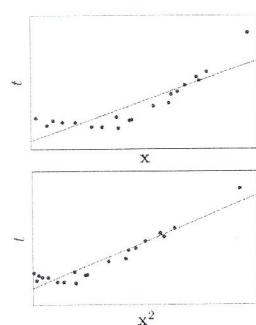
$$y(\mathbf{x}, \mathbf{w}) = w_0 + w_1 \mathbf{x}$$

$\phi_j(\cdot)$  = nonlinear function, called basis function, that is just transformation of the original input.  
We can have how many we want of these basis functions.

Machine Learning

Daniele Lolacano

Let see it in feature space...



For example we can take the green line and add  $x^2$ :  
 $y(\mathbf{x}, \mathbf{w}) = w_0 + w_1 x + w_2 x^2$ .  
In this way:  
 $\int \phi_1(\mathbf{x}) = \mathbf{x}$   
 $\int \phi_2(\mathbf{x}) = \mathbf{x}^2$

The model is still linear

← (in this case should be a plane, in D-dim it should be an hyperplane)

Machine Learning

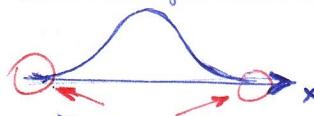
Daniele Lolacano

How can we find basis functions?  
We can do that in a problem-specific way (when we already have some knowledge about the problem and so we know that we could expect some linear correlation between the target and some specific mapping of the input variables) or we can try to apply some standard basis functions:

(the choice of which one to actually use in the model can be done only by comparing the performances of the model)

→ MODEL SELECTION

If we consider a polynomial basis function, e.g.  $\phi(x) = x^2$ , we have that  $\phi(x) \geq 0$  in all the domain. In the case of gaussian we have that if we go far enough from the peak of the gaussian then the function goes almost to zero.

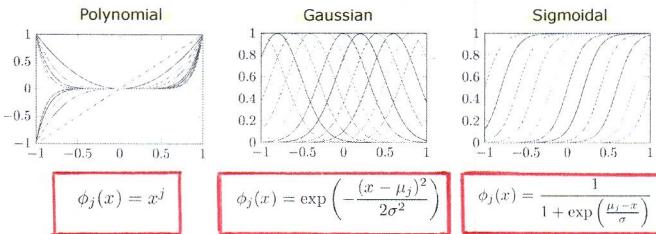


If values are here then their contribution is almost negligible (if values are more extreme then they're far from negligible)  
→ LOCAL (restricted to a neighborhood of the peak of the gaussian)

basis functions are a good way to add previous knowledge in the modelling

### Basis Functions

- Some examples of basis function (assuming single-variable input)

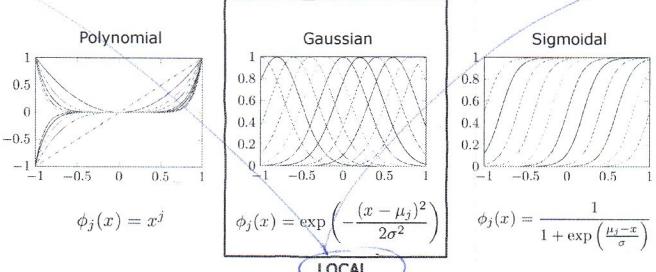


These functions come with a parameter: the grade of the polynomial is a parameter, the mean and variance of the gaussian are parameters,  $\mu_j$  and  $\sigma$  in the sigmoidal are parameters. However, these parameters do not go in  $w$ , they're not learned from data. These parameters are design choices.

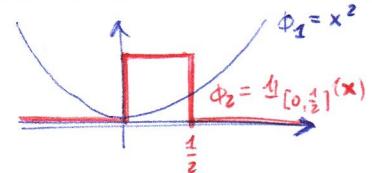
### Machine Learning Daniele Lolicono

### Basis Functions

- Some examples of basis function (assuming single-variable input)

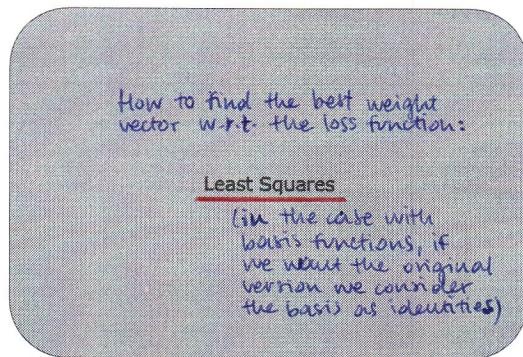


### Machine Learning Daniele Lolicono



In the case  $y(x, w) = w_0 + w_1 \phi_1 + w_2 \phi_2$

No matter what  $x$  is, for sure it will contribute to  $y$  as  $w_1 x^2$ . In:  $y(x, w) = w_0 + w_1 x^2$ , we have that  $x$  will contribute as  $w_1$  only if  $x \in [0, \frac{1}{2}]$ . That's the meaning of LOCAL.



### Machine Learning Daniele Lolicono

### Ordinary Least Squares

- For linear models, a closed-form optimization of the RSS, known as least squares, starting from the matrix form of the loss function:

$$L(w) = \frac{1}{2} RSS(w) = \frac{1}{2}(t - \Phi w)^T(t - \Phi w)$$

where  $\Phi = (\phi(x_1), \dots, \phi(x_N))^T$  and  $t = (t_1, \dots, t_N)^T$

$\varepsilon := (t - \Phi w) = \text{vector of errors}$   
 $(\varepsilon^T \varepsilon = \|\varepsilon\|_2^2)$

$$\Phi = \begin{bmatrix} 1 & \phi_1(x^1) & \phi_2(x^1) & \dots & \phi_{M-1}(x^1) \\ 1 & \phi_1(x^2) & \phi_2(x^2) & \dots & \phi_{M-1}(x^2) \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & \phi_1(x^n) & \phi_2(x^n) & \dots & \phi_{M-1}(x^n) \end{bmatrix} = \begin{bmatrix} 1 & \phi_{11} & \phi_{12} & \dots & \phi_{1(n-1)} \\ 1 & \phi_{21} & \phi_{22} & \dots & \phi_{2(n-1)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & \phi_{n1} & \phi_{n2} & \dots & \phi_{n(n-1)} \end{bmatrix}, \quad t = \begin{bmatrix} t_1 \\ t_2 \\ \vdots \\ t_n \end{bmatrix}, \quad w = \begin{bmatrix} w_0 \\ w_1 \\ \vdots \\ w_{M-1} \end{bmatrix}$$

corresponding target for each sample

vector of weights

### Machine Learning Daniele Lolicono

### Ordinary Least Squares

- For linear models, a closed-form optimization of the RSS, known as least squares, starting from the matrix form of the loss function:

$$L(w) = \frac{1}{2} RSS(w) = \frac{1}{2}(t - \Phi w)^T(t - \Phi w)$$

:  $L(\cdot)$  is quadratic w.r.t.  $w$

where  $\Phi = (\phi(x_1), \dots, \phi(x_N))^T$  and  $t = (t_1, \dots, t_N)^T$

$$\frac{\partial (f(x))^2}{\partial x} = 2 f(x) \cdot f'(x)$$

- We can compute first a second derivative of  $L(w)$  to find the optimal  $w$

$$\frac{\partial L(w)}{\partial w} = -\Phi^T(t - \Phi w)$$

$$\frac{\partial^2 L(w)}{\partial w \partial w^T} = \Phi^T \Phi$$

The concept of convexity is achieved only when this matrix is positive semi-definite (the eigenvalues are  $\geq 0$ ). A matrix of the form  $\Phi^T \Phi$  is for sure positive semi-definite.

since we're in a  $\mathbb{H}$  situation, the stationary point will be a minimum

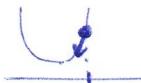
### Machine Learning Daniele Lolicono

### Basis Functions

**Note:**  $\Phi^T \Phi$  is positive semidefinite. In order to invert the matrix we need it to be positive definite  
 $\Rightarrow \hat{w}_{OLS}$  exists when  $\Phi^T \Phi$  is invertible, and so when it is positive definite (all eigenvalues  $> 0$ )

When we have an eigenvalue equal to zero? (not invertible)  
 for example when we have linear dependence among features (we can identify the case and remove the redundant feature).  
 The other case (for eigen = 0) is that we have more features than data.

We work on one sample. We compute the loss function on the one sample we're working on. We compute the derivative of the loss function w.r.t. only that one sample and we apply a correction:



\* We want the learning rate to not decrease too fast (because if so we don't have enough time to reach the optimal solution) but we want it to decrease fast enough (otherwise we don't have convergence).

Ordinary Least Squares

- For linear models, a closed-form optimization of the RSS, known as **least squares**, starting from the matrix form of the loss function:

$$L(\mathbf{w}) = \frac{1}{2} RSS(\mathbf{w}) = \frac{1}{2} (\mathbf{t} - \Phi \mathbf{w})^T (\mathbf{t} - \Phi \mathbf{w})$$

► where  $\Phi = (\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_N))^T$  and  $\mathbf{t} = (t_1, \dots, t_N)^T$

We can compute first a second derivative of  $L(\mathbf{w})$  to find the optimal  $\mathbf{w}$

$$\frac{\partial L(\mathbf{w})}{\partial \mathbf{w}} = -\Phi^T (\mathbf{t} - \Phi \mathbf{w}) \quad \frac{\partial^2 L(\mathbf{w})}{\partial \mathbf{w} \partial \mathbf{w}^T} = \Phi^T \Phi$$

$$\Rightarrow \hat{\mathbf{w}}_{OLS} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$

Assuming  $(\Phi^T \Phi)$  non singular  
Complexity  $O(NM^2 + M^3)$

It's very expensive

Machine Learning

Daniela Loloceno

Sequential Learning / Online learning : when we don't have all the data in advance and we want to update the training each time we get additional data or when we cannot perform OLS because it's not feasible

- Closed-form optimization (OLS) is not feasible with large dataset
- Instead, a stochastic (or sequential) gradient descent is possible
- Least Mean Square (LMS) algorithm:

$$L(\mathbf{x}) = \sum_n L(x_n)$$

$$\Rightarrow \mathbf{w}^{(n+1)} = \mathbf{w}^{(n)} - \alpha^{(n)} \nabla L(x_n)$$

$$\Rightarrow \mathbf{w}^{(n+1)} = \mathbf{w}^{(n)} - \alpha^{(n)} (\mathbf{w}^{(n)T} \phi(\mathbf{x}_n) - t_n) \phi(\mathbf{x}_n)$$

$\alpha$  is called learning rate and to guarantee convergence:

$$\sum_{n=0}^{\infty} \alpha^{(n)} = +\infty \quad \sum_{n=0}^{\infty} \alpha^{(n)2} < +\infty$$

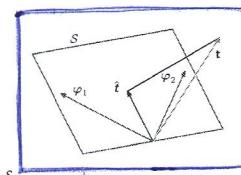
can be variable (admissible) we expect it to decrease

Machine Learning

Daniela Loloceno

### Geometric Interpretation of OLS

- Let  $\mathbf{t}$  be the N-dimensional target vector
- Let  $\varphi_j$  be the  $j$ -th column of matrix  $\Phi$   
 $\varphi_1, \dots, \varphi_M$  identify a linear subspace  $S$
- Let  $\hat{\mathbf{t}}$  be the N-dimensional vector computed as  $\Phi \mathbf{w}$   
 $\hat{\mathbf{t}}$  is a linear combination of  $\varphi_j$  and lies in  $S$
- OLS finds  $\hat{\mathbf{t}}$  minimizing the SSE with respect to  $t$   
 $\hat{\mathbf{t}}$  represents the projection of  $\mathbf{t}$  onto the subspace  $S$



output of our model:  
 $\hat{\mathbf{t}} = \mathbf{y}(\Phi, \hat{\mathbf{w}}_{OLS})$

$$\hat{\mathbf{t}} = \Phi \hat{\mathbf{w}} = \Phi (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$

Hat Matrix (H)

this means that  $\hat{\mathbf{t}}$  is just a linear combination of the columns of  $\Phi$

Daniela Loloceno

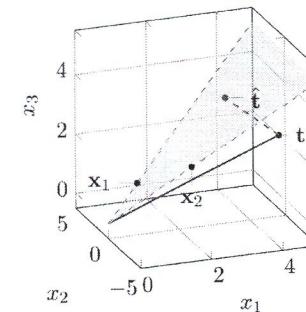
### Geometric Interpretation of OLS: an example

- Let  $N=3$  and  $M=2$

$$\Phi = X = \begin{pmatrix} \varphi_1 & \varphi_2 \\ 1 & 2 \\ 1 & -2 \\ 1 & 2 \end{pmatrix}$$

$$\mathbf{t} = \begin{pmatrix} 5 \\ 1 \\ 2 \end{pmatrix} \quad \text{real target}$$

$$\hat{\mathbf{t}} = \begin{pmatrix} 3.5 \\ 1 \\ 3.5 \end{pmatrix} \quad \text{approx.}$$



Machine Learning

Daniela Loloceno

Multiple Outputs → we want to predict multiple values

- What happens if our regression problem has multiple outputs, i.e.,  $\mathbf{t}$  is not scalar
- It is possible to solve independently a regression problem for each problem
- Yet, it is possible to use the same set of basis functions:

$$\hat{\mathbf{W}} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{T}$$

$$\mathbf{T} = \begin{bmatrix} t_{11} & t_{12} & \dots & t_{1m} \\ t_{21} & t_{22} & \dots & t_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ t_{n1} & t_{n2} & \dots & t_{nm} \end{bmatrix}$$

m-th target of the first sample

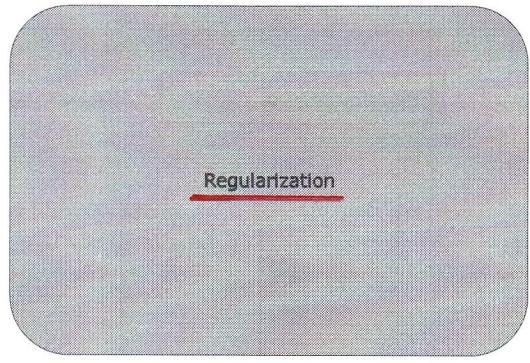
m-th target of the n-th sample

- Where each column of matrix  $\mathbf{T}$  and  $\hat{\mathbf{W}}$  are respectively the target vector of each and the weight vector for each output

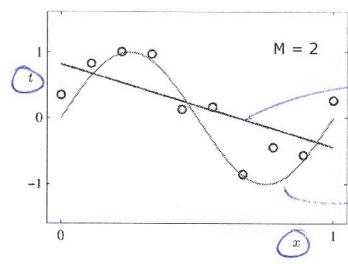
- The solution above can be easily decoupled for each output  $k$ : ( $t_k$  is a vector)

$$\hat{\mathbf{w}}_k = (\Phi^T \Phi)^{-1} \Phi^T t_k$$

► as a benefit  $(\Phi^T \Phi)^{-1}$  can be computed only once



How do we design the linear models? An example

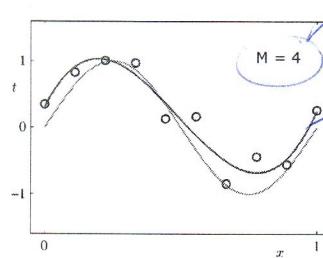


approximation of  $g(\cdot)$ :  
 $y = w_0 + w_1 x$

real function  
from which we  
have the samples

$t = g(x) + \epsilon$ :  $t$  = targets (samples)  
 $g(\cdot)$  = true function

How do we design the linear models? An example



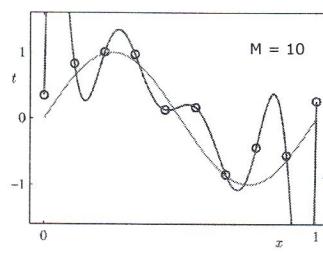
# features  
(including the constant)

approximation obtained  
with least squares using  
a 3-order polynomial model

$$y = w_0 + w_1 x + w_2 x^2 + w_3 x^3$$

3-order polynomial model

How do we design the linear models? An example



With this model the precision  
is accurate only in the already  
given samples → not good  
How can we mitigate? REGULARIZATION

9-order polynomial model

What is regularization?

- What happens to model parameters when complexity increases?

	$M = 1$	$M = 2$	$M = 3$	$M = 10$
$\hat{w}_0$	0.19	0.82	0.31	0.35
$\hat{w}_1$		-1.27	7.99	232.37
$\hat{w}_2$			-25.43	-5321.83
$\hat{w}_3$				48568.31
$\hat{w}_4$				-231639.30
$\hat{w}_5$				640042.26
$\hat{w}_6$				-1061800.52
$\hat{w}_7$				1042400.18
$\hat{w}_8$				-557682.99
$\hat{w}_9$				125201.43

As we add parameters we can see  
that their value increases (in absolute  
value). This is because to have a  
highly-changing function we have  
to have something that changes  
a lot w.r.t.  $x$  → highly-valued  
weights (big values in absolute  
value). To keep the function  
smooth (and not so "crappy") we  
can add an additional term to  
the loss function that has the role  
of keeping the weights as small  
as possible.

## Regularization

- How do we extend the loss function?

$$L(\mathbf{w}) = L_D(\mathbf{w}) + \lambda L_W(\mathbf{w})$$

this additional part has the role of keeping the weights small

- $L_D(\mathbf{w})$  is the usual loss function (e.g., RSS)
- $L_W(\mathbf{w})$  accounts for model complexity
- $\lambda$  is the regularization coefficient

→ the higher the coefficient the more important will be to keep the weights small  
Again: why  $w$  should be small?

If  $w$  is small then the small variations in the feature space will correspond small variations in the target space

This is not a parameter.  
This is a HYPERPARAMETER

- How do we design  $L_w(w)$ ?

- Ridge Regression
- Lasso

Machine Learning

Daniele Lolli cono

## Ridge Regression

- In ridge regression:

$$L_W(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} = \frac{1}{2} \|\mathbf{w}\|_2^2$$

$$\Rightarrow L(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N (t_i - \mathbf{w}^T \phi(x_i))^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2 \quad \|\mathbf{w}\|_2^2 = \sum_{j=0}^{M-1} w_j^2$$

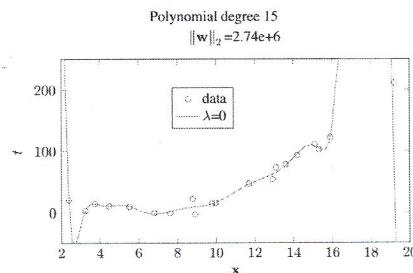
- The loss function is still quadratic with respect to  $w$  and closed-form optimization is still possible:

$$\mathbf{w}_{ridge} = (\lambda \mathbf{I} + \Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$

Machine Learning

Daniele Lolli cono

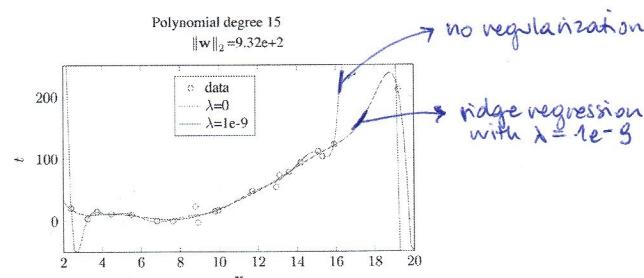
## Ridge Regression: quadratic example



Machine Learning

Daniele Lolli cono

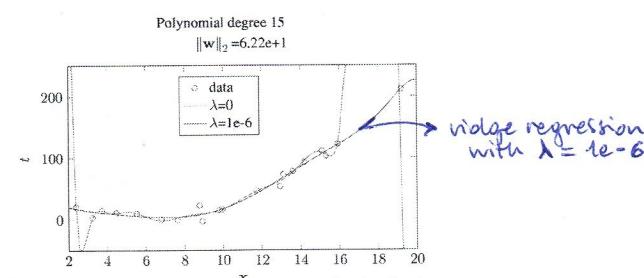
## Ridge Regression: quadratic example



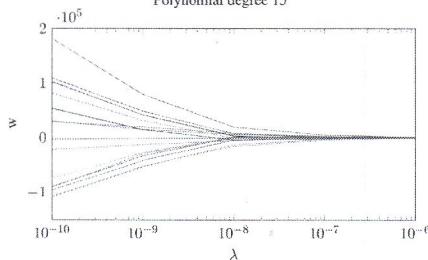
Machine Learning

Daniele Lolli cono

## Ridge Regression: quadratic example



Ridge Regression: quadratic example  
 What happens with the increasing of  $\lambda$ ?  
 The weights decrease in absolute value.  
 Polynomial degree 15



## Machine Learning Daniele Lolacano Lasso

- Another common regularization method is **lasso**:

$$L_W(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|_2^2 = \frac{1}{2} \sum_{j=0}^{M-1} |w_j|^2$$

$$\Rightarrow L(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N (t_i - \mathbf{w}^T \phi(\mathbf{x}_i))^2 + \frac{\lambda}{2} \|\mathbf{w}\|_1 \quad : \quad \|\mathbf{w}\|_1 = \sum_{j=0}^{M-1} |w_j|$$

- In this case, closed-form optimization is not possible
- Nevertheless, lasso typically leads to **sparse** regression models: when regularization coefficient ( $\lambda$ ) is large enough, some **components** of  $\mathbf{w}$  become equal to zero

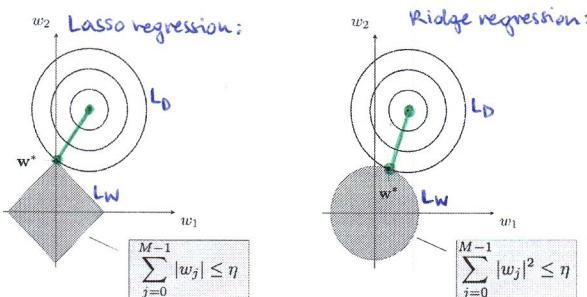
This is very good.  
 If we can have less components then we have simpler models and easier to understand

- We can see regularization equivalent to minimizing  $L_D(w)$  subject to constraint:

$$\sum_{j=0}^{M-1} |w_j| \leq \eta$$

## Machine Learning Daniele Lolacano Lasso vs Ridge Regression

- Why lasso leads to **sparse** model? Let's visualize constraint of lasso and ridge



## Machine Learning Daniele Lolacano Least Squares and Maximum Likelihood (probabilistic approach)

Until now we considered a direct approach. What if we want to proceed in a probabilistic way?

### Least Squares and Maximum Likelihood (probabilistic approach)

## Machine Learning Daniele Lolacano Maximum Likelihood (ML)

- We can deal with regression in a **probabilistic way**
  - We define a probabilistic model that maps inputs ( $x$ ) to outputs ( $t$ )
  - Such probabilistic model,  $y(x, w)$ , will include some **unknown parameters** ( $w$ )
  - Then, we model the **likelihood**, i.e., the probability that observed data  $\mathcal{D}$  is generated by a given set of **parameters** ( $w$ ): *(under the chosen probabilistic model)*

$$p(\mathcal{D}|w)$$

- Finally, we can estimate **parameters** ( $w$ ) by maximizing the likelihood:

$$\mathbf{w}_{ML} = \arg \max_w p(\mathcal{D}|\mathbf{w})$$

## Machine Learning Daniele Lolacano

## Maximum Likelihood (ML) for linear regression

- Our probabilistic model can be defined as:

$$t = y(\mathbf{x}, \mathbf{w}) + \epsilon = \mathbf{w}^T \phi(\mathbf{x}) + \epsilon$$

Given these assumptions and the data we observed, we can evaluate the weights  $\mathbf{w}$ . How? We choose the weights  $\mathbf{w}$  that maximize the probability of the data.

- Given a dataset  $\mathcal{D}$  of  $N$  samples with inputs  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  and outputs  $\mathbf{t} = \{t_1, \dots, t_N\}^\top$ :

$$p(\mathcal{D}|\mathbf{w}) = p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2) = \prod_{n=1}^N \mathcal{N}(t_n | \mathbf{w}^T \phi(\mathbf{x}_n), \sigma^2)$$

known                                  unknown

+ **Hp.** all the points are independent among each other (iid)

goal: maximize this probability w.r.t.  $\mathbf{w}$   
We want to find:

$$\mathbf{w}_{ML} = \arg \max_{\mathbf{w}} p(\mathcal{D}|\mathbf{w})$$

How?  
Instead of maximizing  $p(\mathcal{D}|\mathbf{w})$  we maximize its logarithm.

### Machine Learning Maximum Likelihood (ML) for linear regression (cont.)

Daniele Lolacone

- To find  $\mathbf{w}_{ML}$  it is convenient to maximize the log-likelihood:

$$\mathcal{N}(t_n | \mathbf{w}^T \phi(\mathbf{x}_n), \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp \left\{ -\frac{1}{2\sigma^2} (t_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2 \right\}$$

$$\ell(\mathbf{w}) = \ln p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2) = \sum_{n=1}^N \ln p(t_n | \mathbf{x}_n, \mathbf{w}, \sigma^2) = -\frac{N}{2} \ln (2\pi\sigma^2) - \frac{1}{2\sigma^2} \text{RSS}(\mathbf{w})$$

residual sum of squares

- To solve the optimization problem, we equal the gradient to zero:

$$\nabla \ell(\mathbf{w}) = \sum_{n=1}^N t_n \phi(\mathbf{x}_n)^T - \mathbf{w}^T \left( \sum_{n=1}^N \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^T \right) = 0$$

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



This means that when we apply Ordinary least squares we implicitly do the maximum likelihood approach (with the assumptions of linearity and  $E \sim N(0, \sigma^2)$ )

Daniele Lolacone

This does not provide any tool to get an estimate of the uncertainty level that we have on the estimates

## Bayesian Linear Regression

Machine Learning

Daniele Lolacone

### Bayesian approach

- We formulate our knowledge about the world in a **probabilistic way**:
  - We define the **model** that expresses our knowledge qualitatively
  - Our model will have some **unknown parameters**
  - We capture our **assumptions** about unknown parameters with the **prior distribution** over those parameters before seeing the data
- We observe the **data**
- We compute the posterior probability distribution for the parameters, given observed data

$$p(\text{parameters} | \text{data}) = \frac{p(\text{data} | \text{parameters}) p(\text{parameters})}{p(\text{data})}$$

prior on the parameters  
(before we observe the data what do we know?)

This allows not only to find the most probable value of weights but also to assess the uncertainty that we have in our estimate of the parameters.

### Parameters Posterior Distribution

Daniele Lolacone

- The **posterior distribution** for the model parameters can be found by combining the prior with the likelihood for the parameters given data:

$$p(\mathbf{w} | \mathcal{D}) = \frac{p(\mathcal{D} | \mathbf{w}) p(\mathbf{w})}{p(\mathcal{D})}$$

(Bayes theorem)

overall likelihood of the data. We don't model it, we obtain it by marginalizing  $p(\text{data} | \text{parameters})$

## Parameters Posterior Distribution

- The **posterior distribution** for the model parameters can be found by combining the prior with the likelihood for the parameters given data:

$$p(\mathbf{w}|\mathcal{D}) = \frac{p(\mathcal{D}|\mathbf{w})p(\mathbf{w})}{p(\mathcal{D})}$$

►  $p(\mathbf{w})$  is the **prior probability over the parameter** – what we know before observing the data

## Machine Learning Daniela Lolicono

### Parameters Posterior Distribution

- The **posterior distribution** for the model parameters can be found by combining the prior with the likelihood for the parameters given data:

$$p(\mathbf{w}|\mathcal{D}) = \frac{p(\mathcal{D}|\mathbf{w})p(\mathbf{w})}{p(\mathcal{D})}$$

►  $p(\mathcal{D}|\mathbf{w})$  is the **likelihood** – the probability of observing the data ( $\mathcal{D}$ ) given some value of the parameters ( $\mathbf{w}$ )

## Machine Learning Daniela Lolicono

### Parameters Posterior Distribution

- The **posterior distribution** for the model parameters can be found by combining the prior with the likelihood for the parameters given data:

$$p(\mathbf{w}|\mathcal{D}) = \frac{p(\mathcal{D}|\mathbf{w})p(\mathbf{w})}{p(\mathcal{D})}$$

$$p(\mathcal{D}) = \int p(\mathcal{D}|\mathbf{w})p(\mathbf{w})d\mathbf{w}$$

►  $P(\mathcal{D})$  is the **marginal likelihood** and acts as **normalizing constant**

## Machine Learning Daniela Lolicono

### Parameters Posterior Distribution

- The **posterior distribution** for the model parameters can be found by combining the prior with the likelihood for the parameters given data:

$$p(\mathbf{w}|\mathcal{D}) = \frac{p(\mathcal{D}|\mathbf{w})p(\mathbf{w})}{p(\mathcal{D})}$$

posterior

►  $p(\mathbf{w}|\mathcal{D})$  is the **posterior probability** of parameters  $\mathbf{w}$  given training data  
► the most probable value of  $\mathbf{w}$  given the data will be the mode of the posterior, also known as **maximum a posteriori (MAP)**.

## Machine Learning Daniela Lolicono

### Bayesian Linear Regression

- How to model the **prior**? Assuming a Gaussian likelihood, a **conjugate prior** is the most convenient choice:

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{w}_0, \mathbf{S}_0)$$

→ when we choose a prior we have to choose a mean and a variance for the vector of weights.

How to choose them? it depends on the situation. Often this choice is based on previous knowledge. This settings allow easily to include in the model any previous knowledge we have on the problem. However we can choose what we want for  $\mathbf{w}_0$  and  $\mathbf{S}_0$ .

- As a result, the **posterior** is still Gaussian:

$$p(\mathbf{w}|\mathbf{t}, \Phi, \sigma^2) \propto \mathcal{N}(\mathbf{w}|\mathbf{w}_0, \mathbf{S}_0) \mathcal{N}(\mathbf{t}|\Phi\mathbf{w}, \sigma^2\mathbf{I})$$

(MOST NEUTRAL CHOICE :)

If we don't have any previous knowledge it's good to choose  $\mathbf{w}_0 = \mathbf{0}$  and  $\mathbf{S}_0 = [\sigma^2 \mathbf{0} \mathbf{0} \sigma^2] = \sigma^2 \mathbf{I}$

This is also the choice that we do for the likelihood  
 $(\mathcal{N}(\mathbf{t}|\Phi\mathbf{w}, \sigma^2\mathbf{I}))$

no dependencies between the weights and same variance for all the weights

## Bayesian Linear Regression

How to represent no previous knowledge?  $w_0 = 0$  and  $S_0 = \text{infinity}$  (so we're saying we are infinitely not sure about what we're saying). In this way the posterior mean and variance will depend only on the data:

$$\begin{cases} w_N = S_N \cdot \frac{\Phi^T t}{\sigma^2} \\ S_N^{-1} = \frac{\Phi^T \Phi}{\sigma^2} \end{cases}$$

which, if we substitute, brings:

$$w_N = (\Phi^T \Phi)^{-1} \Phi^T t$$

which is again OLS!

So, maximum likelihood / OLS is the specific case of the Bayesian regression framework that corresponds to the case where we don't have any previous knowledge on weights.

Most easiest and UNBIASED estimate that we can get from data

(notice that we have  $N-M$  because with  $M$  we account the degrees of freedom)

- How to model the prior? Assuming a Gaussian likelihood, a conjugate prior is the most convenient choice:

$$p(w) = \mathcal{N}(w|w_0, S_0)$$

- As a result, the posterior is still Gaussian:

$$p(w|t, \Phi, \sigma^2) \propto \mathcal{N}(w|w_0, S_0) \mathcal{N}(t|\Phi w, \sigma^2 I)$$

$$\rightarrow \begin{cases} p(w|t, \Phi, \sigma^2) = \mathcal{N}(w|w_N, S_N) \\ w_N = S_N \left( S_0^{-1} w_0 + \frac{\Phi^T t}{\sigma^2} \right) \\ S_N^{-1} = S_0^{-1} + \frac{\Phi^T \Phi}{\sigma^2} \end{cases}$$

conjugate  $\rightarrow$  still gaussian

Note: if we have a strong assumption in the prior, this strong assumption will be reflected also in the posterior

## Machine Learning

Daniele Lolocorno

### Bayesian Linear Regression an Maximum Likelihood

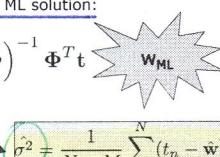
- Which parameters?

- The maximum a-posteriori (MAP) is the obvious choice
- When posterior is gaussian the MAP is equal to the mean

- When prior is infinitely broad MAP is equal to ML solution:

$$\lim_{S_0 \rightarrow \infty} w_N = (\Phi^T \Phi)^{-1} \Phi^T t$$

$$\lim_{S_0 \rightarrow \infty} S_N^{-1} = \frac{\Phi^T \Phi}{\sigma^2}$$



- The ML estimate of  $w$  has the **smallest variance** among linear unbiased estimates and the **lowest MSE** among linear unbiased estimates (Gauss-Markov).

(Theorem)

## Machine Learning

Daniele Lolocorno

### Bayesian Linear Regression and Regularization

- What about regularization?

- When  $w_0=0$  and  $S_0=\tau^2 I$ , we have:

$$\ln p(w|t) = -\frac{1}{2\sigma^2} \sum_{i=1}^N (t_i - w^T \phi(x_i))^2 - \frac{1}{2\tau^2} \|w\|_2^2$$

- In this case, MAP ( $w_N$ ) is equivalent to the solution of ridge regression ( $\hat{w}_{\text{ridge}}$ ) with  $\lambda = \frac{\sigma^2}{\tau^2}$

The great thing about Bayesian approach is that, unlike in the case of OLS or maximum likelihood, here we get two estimates:  $w_N, S_N$ . We always got only  $w_N$ . Now we have an estimate of THE UNCERTAINTY we have on the weight vector.

## Machine Learning

Daniele Lolocorno

### Bayesian Linear Regression: an example

- How to exploit the Bayesian approach for sequential learning?

- We compute **posterior** with initial data
- When additional data is available, the **posterior becomes the prior**

$$\begin{aligned} p(w|D) &= N(w_N, S_N) \\ \text{then we obtain new data: } D' &: \\ p(w'|D') &= p(D'|w') p(w') \end{aligned}$$

The maximum a-posterior ( $w_N$ ) is the same as the result obtained with ridge regression if we assume the prior to have 0 mean and covariance matrix  $= \tau^2 I$ . In this case it'll result:

$$\lambda = \sigma^2 / \tau^2.$$

$\Rightarrow$  Bayesian regression also allows to account for regularization. Regularization can be seen as specific case of Bayesian regression for which we make specific decisions for the prior.

MOREOVER

Bayesian regression can be used "incrementally".

## Machine Learning

Daniele Lolocorno

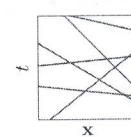
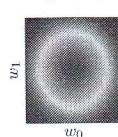
### Bayesian Linear Regression: an example

- How to exploit the Bayesian approach for sequential learning?

- We compute **posterior** with initial data
- When additional data is available, the **posterior becomes the prior**

- Let see an example

- Let assume data is generated as  $t(x) = -0.3 + 0.5x + \varepsilon$
- Let assume that  $x \sim U(-1, 1)$  and  $\varepsilon \sim \mathcal{N}(0, 0.04)$
- Let use as model:  $y(x, w) = w_0 + w_1 x$
- Let assume as prior:  $p(w) = \mathcal{N}(w_0, \tau^2 I)$  with  $\tau^2 = 0.5$  and  $w_0 = [0, 0]^T$

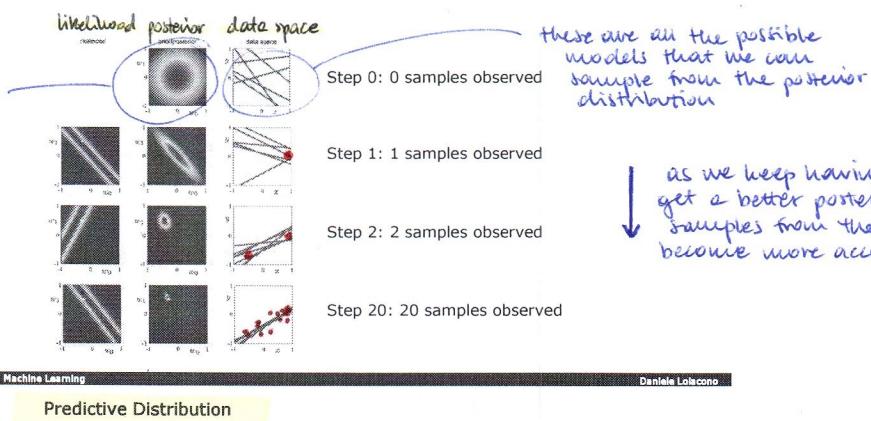


$\rightarrow$  this means that we're basically doing ridge regression

### Bayesian Linear Regression: an example (2)

At the first step the posterior is basically the prior because we don't have any samples

(\*) Note: likelihood and posterior have the weights  $w_0$  and  $w_1$  on the axis. We sample a couple from the joint posterior distribution of  $w_0$  and  $w_1$ , we obtain a sample and so a line. The line is represented in the data space



- With a Bayesian framework, it is possible to compute the probability distribution of the target for a new sample  $x^*$  (given the training data  $\mathcal{D}$ ), by integrating over the posterior distribution:

$$p(t^*|x^*, \mathcal{D}) = \mathbb{E}[t^*|x^*, \mathbf{w}, \mathcal{D}] = \underbrace{\int p(t^*|\mathbf{x}^*, \mathbf{w}, \mathcal{D})p(\mathbf{w}|\mathcal{D})d\mathbf{w}}_{\text{MAP}}$$

- This is often called **predictive distribution**
- Unfortunately, computing the predictive distribution requires knowledge of the posterior distribution, which is usually intractable

However, in our case (linear bayesian regression, gaussian distribution for the prior, gaussian likelihood, etc.) we can obtain a solution analytically:

### Predictive Distribution for Bayesian Regression

(Model assumptions:  $y(x) = \mathbf{w}^T \phi(x) + \epsilon$ ,  $\epsilon \sim N(0, \sigma^2)$ )

- With our assumptions, we can compute the **predictive distribution**:

$$p(t|x, \mathcal{D}, \sigma^2) = \int p(t|\mathbf{x}, \mathbf{w}, \sigma^2)p(\mathbf{w}|\mathbf{w}_N, \mathbf{S}_N)d\mathbf{w} = \int \mathcal{N}(t|\mathbf{w}^T \phi(x), \sigma_N^2) \mathcal{N}(\mathbf{w}|\mathbf{w}_N, \mathbf{S}_N)d\mathbf{w}$$

$$\boxed{p(t|x, \mathcal{D}, \sigma^2) = \mathcal{N}(t|\mathbf{w}_N^T \phi(x), \sigma_N^2(x))}$$

$$\Rightarrow \boxed{\sigma_N^2(x) = \sigma^2 + \phi(x)^T \mathbf{S}_N \phi(x)}$$

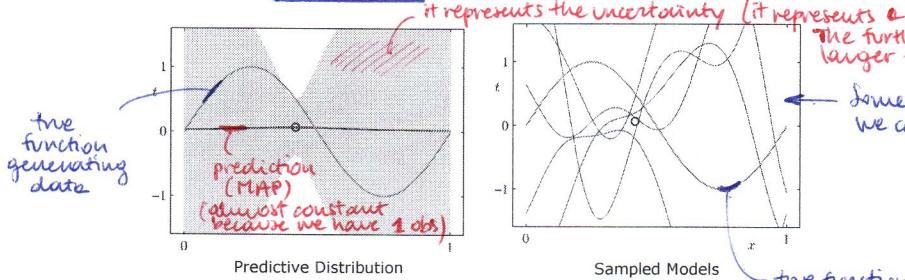
data parameters

- when  $N \rightarrow \infty$  the uncertainty associated to parameters (second term) goes to zero and the variance of predictive distribution depends only on variance of data ( $\sigma^2$ ) (we cannot reduce  $\sigma^2$ : we cannot have a degree of uncertainty on the target that is less than the actual variance of the data)

the mean does not change (it's still the MAP), however we're interested in  $\sigma_N^2(x)$  because that is the one that will tell about variance (thanks to that we can construct a confidence interval!)

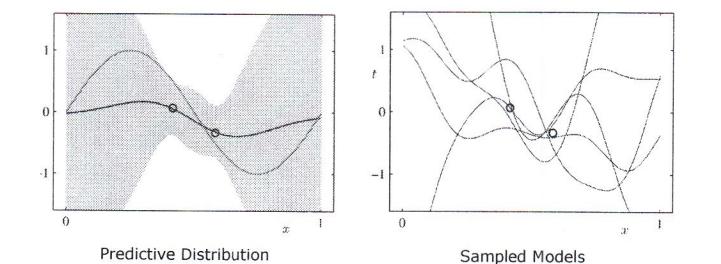
### Predictive Distribution: an example

- Approximating a sinusoidal dataset with a linear model w/ 9 Gaussian basis functions: 1 sample observed



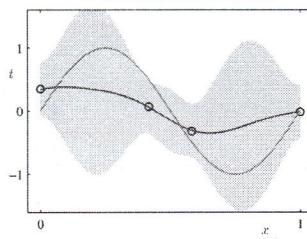
### Predictive Distribution: an example

- Approximating a sinusoidal dataset with a linear model w/ 9 Gaussian basis functions: 2 samples observed

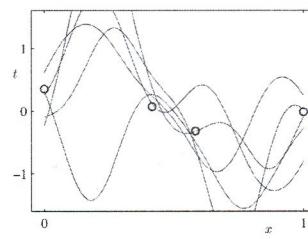


## Predictive Distribution: an example

- Approximating a sinusoidal dataset with a linear model w/ 9 Gaussian basis functions: 4 samples observed



Predictive Distribution



Sampled Models

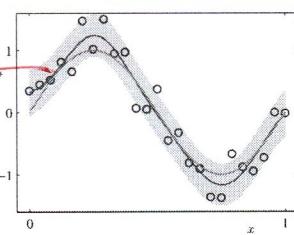
Machine Learning

Daniele Lalaconi

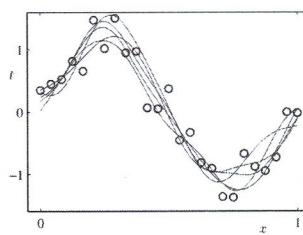
## Predictive Distribution: an example

- Approximating a sinusoidal dataset with a linear model w/ 9 Gaussian basis functions: 25 samples observed

*Much smaller confidence intervals!*



Predictive Distribution



Sampled Models

Machine Learning

Daniele Lalaconi

## Challenges and Limitation of Linear Regression

Machine Learning

Daniele Lalaconi

### Challenges

- Modeling Challenges
  - Our model should fit well all the functions (we think) are **likely**
  - The prior should not give zero or small probabilities to possible values, but at the same time also avoid spreading out the probability (**uninformative**)
- Computational Challenges
  - **Analytical Integration** is possible only using **conjugate priors** and works for **simple models**
  - **Approximated** approaches are instead to be used in the more general case, such as Gaussian (Laplace) approximation, Monte Carlo integration, and Variational approximation

Machine Learning

Daniele Lalaconi

### Limitation of Fixed Basis Functions

- Linear models with fixed basis functions have several advantages
  - Allow closed-form solution
  - Tractable **Bayesian treatment**
  - Can model non-linear relationship with proper basis function
- However, they have also several limitations
  - Basis functions are not *adaptive* with respect to the training data
  - These models suffer of the **curse of dimensionality**

**BIGGEST INCONVENIENT:**  
(we cannot adapt  $\phi(\cdot)$  to the data,  
we have TO choose all of them in  
advance)

Machine Learning

Daniele Lalaconi