

# Betting on Sparsity with the Lasso. Part II: More on ridge/lasso and beyond linear regression



Pierre-Alexandre Mattei

<http://pamattei.github.io> – @pamattei  
pierre-alexandre.mattei@inria.fr

MSc Data Science

Recap on regularised linear regression

Basic lasso theory

Recap on regularised linear regression

Basic lasso theory

## Good old (centered) linear regression

---

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

If we have a new data point  $\mathbf{x}_{\text{new}}$ , how do we predict its response?

## Good old (centered) linear regression

---

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

If we have a new data point  $\mathbf{x}_{\text{new}}$ , how do we predict its response?

We can just use

$$\hat{\mathbf{y}}_{\text{new}} = \mathbb{E}[y|\mathbf{x}_{\text{new}}] = \mathbf{x}_{\text{new}}^T \boldsymbol{\beta}.$$

## Good old (centered) linear regression

---

$$\mathbf{Y} = \mathbf{X}\beta + \epsilon$$

If we have a new data point  $\mathbf{x}_{\text{new}}$ , how do we predict its response?

We can just use

$$\hat{\mathbf{y}}_{\text{new}} = \mathbb{E}[y|\mathbf{x}_{\text{new}}] = \mathbf{x}_{\text{new}}^T \beta.$$

If we mainly want to do prediction, we're mostly interested in estimating  $\beta$  (and the noise term is not too important to estimate).

# OLS for linear regression

---

Since Legendre and Gauss ( $\approx 1805$ ), the traditional approach to linear regression is through **ordinary least squares**

$$\hat{\beta}_{\text{OLS}} \in \operatorname{argmin}_{\beta \in \mathbb{R}^p} \|\mathbf{Y} - \mathbf{X}\beta\|_2^2,$$

which may be interpreted as doing maximum likelihood under the assumption that the noise term is Gaussian.

# OLS for linear regression

---

Since Legendre and Gauss ( $\approx 1805$ ), the traditional approach to linear regression is through **ordinary least squares**

$$\hat{\beta}_{\text{OLS}} \in \operatorname{argmin}_{\beta \in \mathbb{R}^p} \|\mathbf{Y} - \mathbf{X}\beta\|_2^2,$$

which may be interpreted as doing maximum likelihood under the assumption that the noise term is Gaussian. If the matrix  $\mathbf{X}^T \mathbf{X}$  is invertible, then the problem admits a **unique solution**:

$$\hat{\beta}_{\text{OLS}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}.$$



# OLS for linear regression

---

Since Legendre and Gauss ( $\approx 1805$ ), the traditional approach to linear regression is through **ordinary least squares**

$$\hat{\beta}_{\text{OLS}} \in \operatorname{argmin}_{\beta \in \mathbb{R}^p} \|\mathbf{Y} - \mathbf{X}\beta\|_2^2,$$

which may be interpreted as doing maximum likelihood under the assumption that the noise term is Gaussian. If the matrix  $\mathbf{X}^T \mathbf{X}$  is invertible, then the problem admits a **unique solution**:

$$\hat{\beta}_{\text{OLS}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}.$$

In this course, we're concerned by **cases where the number of features  $p$  is very large**. Which leads  $\mathbf{X}^T \mathbf{X}$  to be ill-conditioned or non-invertible. **This renders OLS impractical.**

## OLS fails in high dimensions

---

We have  $\text{rank}(\mathbf{X}^T \mathbf{X}) \leq n$ , so if  $p > n$ , the matrix  $\mathbf{X}^T \mathbf{X}$  is not invertible!  
There is an infinite number of minimisers of the squared error...

# OLS fails in high dimensions

---

We have  $\text{rank}(\mathbf{X}^T \mathbf{X}) \leq n$ , so if  $p > n$ , the matrix  $\mathbf{X}^T \mathbf{X}$  is not invertible!  
There is an infinite number of minimisers of the squared error...

**Exercise:** We saw last year three simple ways of doing that. What were they?

# OLS fails in high dimensions

---

We have  $\text{rank}(\mathbf{X}^T \mathbf{X}) \leq n$ , so if  $p > n$ , the matrix  $\mathbf{X}^T \mathbf{X}$  is not invertible!  
There is an infinite number of minimisers of the squared error...

**Exercise:** We saw last year three simple ways of doing that. What were they?

- replacing  $(\mathbf{X}^T \mathbf{X})^{-1}$  by the **Moore-Penrose pseudoinverse**  $(\mathbf{X}^T \mathbf{X})^\dagger$   
("ridgeless regression" or "Moore-Penrose least squares")

# OLS fails in high dimensions

---

We have  $\text{rank}(\mathbf{X}^T \mathbf{X}) \leq n$ , so if  $p > n$ , the matrix  $\mathbf{X}^T \mathbf{X}$  is not invertible!  
There is an infinite number of minimisers of the squared error...

**Exercise:** We saw last year three simple ways of doing that. What were they?

- replacing  $(\mathbf{X}^T \mathbf{X})^{-1}$  by the **Moore-Penrose pseudoinverse**  $(\mathbf{X}^T \mathbf{X})^\dagger$  ("ridgeless regression" or "Moore-Penrose least squares")
- replacing  $(\mathbf{X}^T \mathbf{X})^{-1}$  by  $(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p)^{-1}$  with  $\lambda > 0$  ("**ridge regression**", "Tikhonov regularisation", " $\ell_2$  regularisation")

# OLS fails in high dimensions

---

We have  $\text{rank}(\mathbf{X}^T \mathbf{X}) \leq n$ , so if  $p > n$ , the matrix  $\mathbf{X}^T \mathbf{X}$  is not invertible!  
There is an infinite number of minimisers of the squared error...

**Exercise:** We saw last year three simple ways of doing that. What were they?

- replacing  $(\mathbf{X}^T \mathbf{X})^{-1}$  by the **Moore-Penrose pseudoinverse**  $(\mathbf{X}^T \mathbf{X})^\dagger$  ("ridgeless regression" or "Moore-Penrose least squares")
- replacing  $(\mathbf{X}^T \mathbf{X})^{-1}$  by  $(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_p)^{-1}$  with  $\lambda > 0$  ("**ridge regression**", "Tikhonov regularisation", " $\ell_2$  regularisation")
- adding an  $\ell_0$  penalty to the sum of squared errors ("**lasso**", "basis pursuit").

**Exercise:** What are the advantages/drawbacks of the three methods?

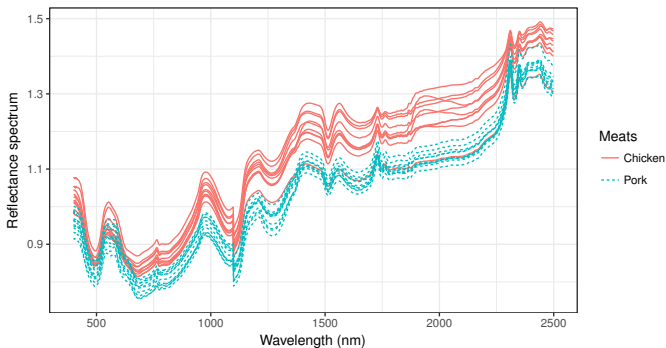
## Betting on sparsity

---

A drawback of both Moore-Penrose and ridge is that they **do not lead to a sparse solution**. In other words, for them, all  $p$  variables are relevant.

# Betting on sparsity

A drawback of both Moore-Penrose and ridge is that they **do not lead to a sparse solution**. In other words, for them, all  $p$  variables are relevant. For example, when dealing with **spectra**, not all wavelengths are useful in general.





# What does sparsity actually mean?

---

A **sparse model** is a model that willingly ignores some of the features of the data at hand. For linear regression, this means that the parameter  $\beta$  will have some coefficients equal to zero.

# What does sparsity actually mean?

---

A **sparse model** is a model that willingly ignores some of the features of the data at hand. For linear regression, this means that the parameter  $\beta$  will have some coefficients equal to zero. It will be convenient to use the  $\ell_0$  pseudo-norm of a vector  $\beta \in \mathbb{R}^p$ , defined as

$$\|\beta\|_0 = \#\{\beta_j \mid \beta_j \neq 0\} = \text{number of nonzero coefficients of } \beta.$$

# What does sparsity actually mean?

---

A **sparse model** is a model that willingly ignores some of the features of the data at hand. For linear regression, this means that the parameter  $\beta$  will have some coefficients equal to zero. It will be convenient to use the  $\ell_0$  pseudo-norm of a vector  $\beta \in \mathbb{R}^p$ , defined as

$$\|\beta\|_0 = \#\{\beta_j \mid \beta_j \neq 0\} = \text{number of nonzero coefficients of } \beta.$$

We say that  $\beta \in \mathbb{R}^p$  is  **$k$ -sparse** when it contains only  $k$  coefficients different from zero. In other words,  $\beta \in \mathbb{R}^p$  is  $k$ -sparse when  $\|\beta\|_0 = k$ . Of course, we need to have  $k \in \{0, \dots, p\}$ . The **sparsity pattern** of  $\beta$  is the subset of  $\{1, \dots, p\}$  corresponding to nonzero coefficients.

# What does sparsity actually mean?

---

A **sparse model** is a model that willingly ignores some of the features of the data at hand. For linear regression, this means that the parameter  $\beta$  will have some coefficients equal to zero. It will be convenient to use the  $\ell_0$  pseudo-norm of a vector  $\beta \in \mathbb{R}^p$ , defined as

$$\|\beta\|_0 = \#\{\beta_j \mid \beta_j \neq 0\} = \text{number of nonzero coefficients of } \beta.$$

We say that  $\beta \in \mathbb{R}^p$  is  **$k$ -sparse** when it contains only  $k$  coefficients different from zero. In other words,  $\beta \in \mathbb{R}^p$  is  $k$ -sparse when  $\|\beta\|_0 = k$ . Of course, we need to have  $k \in \{0, \dots, p\}$ . The **sparsity pattern** of  $\beta$  is the subset of  $\{1, \dots, p\}$  corresponding to nonzero coefficients.

There are  $2^p$  possible sparsity patterns. This means that, for 20 features, we'll already have more than one million patterns... **Trying all of them will be impossible in "large  $p$ " cases!**

## Sparsity through $\ell_0$ penalisation

---

Since we have a measure of non-sparsity (the  $\ell_0$  pseudonorm), we could just use that to obtain sparse solutions!

$$\hat{\beta}_{\ell_0, \lambda} \in \operatorname{argmin}_{\beta \in \mathbb{R}^p} \|\mathbf{Y} - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_0.$$

## Sparsity through $\ell_0$ penalisation

---

Since we have a measure of non-sparsity (the  $\ell_0$  pseudonorm), we could just use that to obtain sparse solutions!

$$\hat{\beta}_{\ell_0, \lambda} \in \operatorname{argmin}_{\beta \in \mathbb{R}^p} \|\mathbf{Y} - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_0.$$

If we use a large enough  $\lambda$ , this will give us sparse solutions. However, the loss function is **non-differentiable**, if we wanted to solve the problem exactly, we would basically need to try all possible  $2^p$  sparsity patterns and compute OLS on them. **This is impossible when  $p$  becomes bigger than around 50...**

Since we can't do  $\ell_0$  regularisation, what can we do?

---

Another way of seeing the  $\ell_0$  regularised problem is as an  $\ell_0$  constrained problem

$$\operatorname{argmin}_{\beta \in \mathbb{R}^p, \|\beta\|_0 \leq k} \|\mathbf{Y} - \mathbf{X}\beta\|_2^2.$$

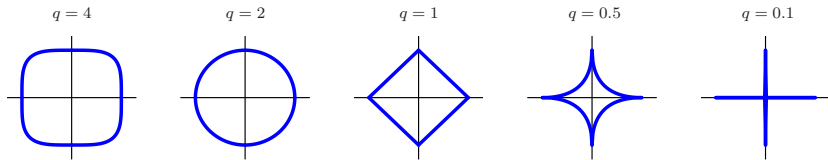
What does the  $\ell_0$  ball  $\{\beta \in \mathbb{R}^p, \|\beta\|_0 \leq k\}$  look like?

## The $\ell_q$ pseudonorms

One neat way of gaining insight on the  $\ell_0$  ball is to see it as a **limit of  $\ell_q$  balls**. For all  $q > 0$ , let us define

$$\|\beta\|_q = \left( \sum_{j=1}^p \beta_j^q \right)^{1/q}.$$

Note that this measure is not a proper norm unless  $q \geq 1$ . The  $\ell_0$  case corresponds to  $q \rightarrow 0$ . Here are various  $\ell_q$  balls (figure from Hastie, Tibshirani & Wainwright, 2015).



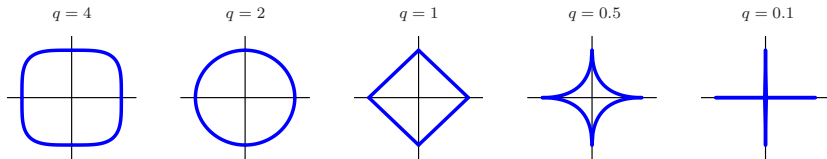


# The $\ell_q$ pseudonorms

One neat way of gaining insight on the  $\ell_0$  ball is to see it as a **limit of  $\ell_q$  balls**. For all  $q > 0$ , let us define

$$\|\beta\|_q = \left( \sum_{j=1}^p \beta_j^q \right)^{1/q}.$$

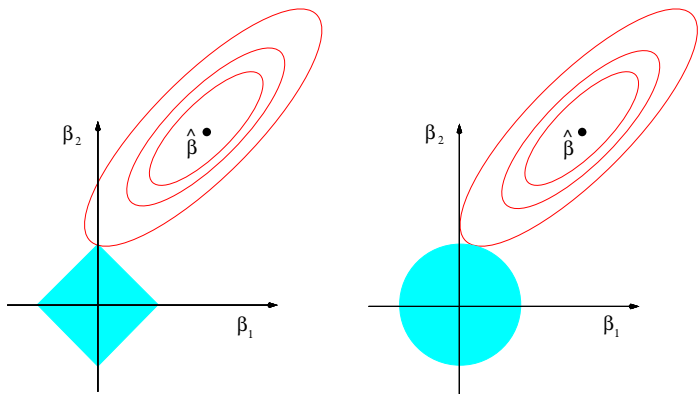
Note that this measure is not a proper norm unless  $q \geq 1$ . The  $\ell_0$  case corresponds to  $q \rightarrow 0$ . Here are various  $\ell_q$  balls (figure from Hastie, Tibshirani & Wainwright, 2015).



**Key idea of the lasso: replace the discrete, non-differentiable, non-convex  $\ell_0$  ball by a more regular object, like an  $\ell_q$  ball.**

## From ridge to lasso

We already studied the  $\ell_2$  case, it's just ridge! But **ridge does not give sparsity....** To get sparsity, we need **sharp edges on the ball** (figure from Hastie, Tibshirani & Wainwright, 2015).



## From ridge to lasso

---

All  $\ell_q$  balls have sharp edges when  $q < 1$ , and all of them would lead to sparse solutions.

# From ridge to lasso

---

All  $\ell_q$  balls have sharp edges when  $q < 1$ , and all of them would lead to sparse solutions.

Keep in mind that we're doing this for **computational reasons**: we want something fast and cheap.

This leads to the desideratum of having a **convex optimisation problem**. Of course, the squared error is convex in  $\beta$ . What about the penalty?

## From ridge to lasso

---

All  $\ell_q$  balls have sharp edges when  $q < 1$ , and all of them would lead to sparse solutions.

Keep in mind that we're doing this for **computational reasons**: we want something fast and cheap.

This leads to the desideratum of having a **convex optimisation problem**. Of course, the squared error is convex in  $\beta$ . What about the penalty?

The only convex  $\ell_q$  ball is the  $\ell_1$  ball, which justifies the choice of  $q = 1$ . This leads to the **lasso estimate**

$$\hat{\beta}_{\text{lasso},\lambda} \in \operatorname{argmin}_{\beta \in \mathbb{R}^p} \|\mathbf{Y} - \mathbf{X}\beta\|_2^2 + \lambda \|\beta\|_1.$$

Recap on regularised linear regression

Basic lasso theory

# What does it mean to do lasso theory?

---

There are many goals one could have in mind. Here are a few examples, on the **fundamental side**...

- **estimation**: assuming the "true model" is really sparse, does the lasso solution recover the true sparsity pattern

# What does it mean to do lasso theory?

---

There are many goals one could have in mind. Here are a few examples, on the **fundamental side...**

- **estimation:** assuming the "true model" is really sparse, does the lasso solution recover the true sparsity pattern
- **support recovery:** assuming the "true model" is really sparse, does the lasso solution recover the true sparsity pattern



# What does it mean to do lasso theory?

---

There are many goals one could have in mind. Here are a few examples, on the **fundamental side...**

- **estimation:** assuming the "true model" is really sparse, does the lasso solution recover the true sparsity pattern
- **support recovery:** assuming the "true model" is really sparse, does the lasso solution recover the true sparsity pattern
- **predictive performance:** does the lasso solution give better predictions than other (linear) predictions?

# What does it mean to do lasso theory?

---

There are many goals one could have in mind. Here are a few examples, on the **fundamental side**...

- **estimation**: assuming the "true model" is really sparse, does the lasso solution recover the true sparsity pattern
- **support recovery**: assuming the "true model" is really sparse, does the lasso solution recover the true sparsity pattern
- **predictive performance**: does the lasso solution give better predictions than other (linear) predictions?

An important question is whether or not one has to make **asymptotic assumptions** (i.e.  $n$ ,  $p$ , or both go to infinity) in order to have theoretical results. **Non-asymptotic results** are usually more realistic and useful, but they are harder to derive and sometimes uglier/more cryptic. All these results can be useful to know **under which assumptions the lasso "works"**, and when it is relevant to practically use it.

# What does it mean to do lasso theory?

---

... and on a more **practical side**...

- **optimisation**: is the lasso solution unique? what are the guarantees that we will find it in reasonable time?
- **degrees of freedom**: can we derive AIC/BIC-type criteria for the lasso?
- **properties of  $\hat{\beta}_{\text{lasso}}$** : are we sure that some coefficients will be zero? how many of them?
- can we use theoretical insight to **improve the lasso**, or design **faster algorithms** for it?

## Brief history of lasso theory

---

Important theoretical work on the lasso has been one the biggest trends in stats of the period 2000-2020.

## Brief history of lasso theory

---

Important theoretical work on the lasso has been one the biggest trends in stats of the period 2000-2020.

One of the first important papers that started it was *Asymptotics for lasso-type estimators*, by Knight & Fu (*Annals of Statistics*, 2000).

Then, many prominent researchers from stats/proba/optim communities tackled the issues of the previous slides, even including "famous outsiders" like Terence Tao (Candès & Tao, *Annals of Statistics*, 2007).

## Brief history of lasso theory

---

Important theoretical work on the lasso has been one the biggest trends in stats of the period 2000-2020.

One of the first important papers that started it was *Asymptotics for lasso-type estimators*, by Knight & Fu (*Annals of Statistics*, 2000).

Then, many prominent researchers from stats/proba/optim communities tackled the issues of the previous slides, even including "famous outsiders" like Terence Tao (Candès & Tao, *Annals of Statistics*, 2007).

A good overview of these results is provided in Chapter 11 of Hastie, Tibshirani & Wainwright (2015).

# The kind of theory we'll study today

---

We will show a simple theorem that allows us to understand where the zeros in the lasso solution come from.

# The kind of theory we'll study today

---

We will show a **simple theorem that allows us to understand where the zeros in the lasso solution come from.**

We will state without proof a few interesting (and much harder to prove) results. Some more fundamental, some more practical.



## Exercise: Tibshirani's simple theorem

---

In the original 1996 lasso paper, Robert Tibshirani stated (without proof) a cute little result that gives valuable insight about why the lasso solution is sparse. That's the subject of this exercise.

## Exercise: Tibshirani's simple theorem

---

In the original 1996 lasso paper, Robert Tibshirani stated (without proof) a cute little result that gives valuable insight about why the lasso solution is sparse. That's the subject of this exercise.

**Assumption (orthogonal design):** we assume that  $\mathbf{X}^T \mathbf{X} = \mathbf{I}_p$ .

## Exercise: Tibshirani's simple theorem

---

In the original 1996 lasso paper, Robert Tibshirani stated (without proof) a cute little result that gives valuable insight about why the lasso solution is sparse. That's the subject of this exercise.

**Assumption (orthogonal design):** we assume that  $\mathbf{X}^T \mathbf{X} = \mathbf{I}_p$ .

1. Discuss this assumption.

## Exercise: Tibshirani's simple theorem

---

In the original 1996 lasso paper, Robert Tibshirani stated (without proof) a cute little result that gives valuable insight about why the lasso solution is sparse. That's the subject of this exercise.

**Assumption (orthogonal design):** we assume that  $\mathbf{X}^T \mathbf{X} = \mathbf{I}_p$ .

1. Discuss this assumption. In particular, is it compatible with  $p > n$ ?

## Exercise: Tibshirani's simple theorem

---

In the original 1996 lasso paper, Robert Tibshirani stated (without proof) a cute little result that gives valuable insight about why the lasso solution is sparse. That's the subject of this exercise.

**Assumption (orthogonal design):** we assume that  $\mathbf{X}^T \mathbf{X} = \mathbf{I}_p$ .

1. Discuss this assumption. In particular, is it compatible with  $p > n$ ?
2. Write down the OLS estimate  $\hat{\beta}_{\text{OLS}}$  in this orthogonal case.
3. Show that the lasso problem is equivalent to a problem of the form

$$\max_{\beta \in \mathbb{R}^p} \sum_{j=1}^p f_j(\beta_j),$$

where the  $f_j$ s are simple (piecewise polynomial) functions. 4. Find the explicit form of  $\hat{\beta}_{\text{lasso}}$ . *Possible hint: you may start by assuming that all coefficients of  $\hat{\beta}_{\text{OLS}}$  are positive.*

# Tibshirani's simple theorem

---

Theorem (Tibshirani, 1996). If  $\mathbf{X}^T \mathbf{X} = \mathbf{I}_p$ , then

$$\hat{\beta}_{\text{lasso}, \lambda} = \text{sign}(\hat{\beta}_{\text{OLS}})(|\hat{\beta}_{\text{OLS}}| - \lambda)^+,$$

where  $x^+ = \max(x, 0) = \text{ReLU}(x)$ .

---

<sup>1</sup>[https://stats.stackexchange.com/questions/323234/  
how-is-the-lasso-orthogonal-design-case-solution-derived](https://stats.stackexchange.com/questions/323234/how-is-the-lasso-orthogonal-design-case-solution-derived)

# Tibshirani's simple theorem

---

Theorem (Tibshirani, 1996). If  $\mathbf{X}^T \mathbf{X} = \mathbf{I}_p$ , then

$$\hat{\beta}_{\text{lasso}, \lambda} = \text{sign}(\hat{\beta}_{\text{OLS}})(|\hat{\beta}_{\text{OLS}}| - \lambda)^+,$$

where  $x^+ = \max(x, 0) = \text{ReLU}(x)$ .

The proof that inspired my exercise comes from stats stackexchange!<sup>1</sup>

---

<sup>1</sup><https://stats.stackexchange.com/questions/323234/how-is-the-lasso-orthogonal-design-case-solution-derived>

# Tibshirani's simple theorem

---

**Theorem (Tibshirani, 1996).** If  $\mathbf{X}^T \mathbf{X} = \mathbf{I}_p$ , then

$$\hat{\beta}_{\text{lasso}, \lambda} = \text{sign}(\hat{\beta}_{\text{OLS}})(|\hat{\beta}_{\text{OLS}}| - \lambda)^+,$$

where  $x^+ = \max(x, 0) = \text{ReLU}(x)$ .

The proof that inspired my exercise comes from stats stackexchange!<sup>1</sup>

The function  $S : x \mapsto \text{sign}(x)(|x| - \lambda)^+$  is called the **soft-thresholding operator**, and is present a lot in sparse contexts.

---

<sup>1</sup><https://stats.stackexchange.com/questions/323234/how-is-the-lasso-orthogonal-design-case-solution-derived>



## A few more complex results

---

There are a lot of result that show that, under some (unfortunately strong) assumptions on  $\mathbf{X}$ , we can have consistency of  $\hat{\beta}_{\text{lasso},\lambda}$  or of its support when both  $n$  and  $p$  are large but the model is assumed to be truly sparse. See chapter 11 of of Hastie, Tibshirani & Wainwright (2015) for more on this.

## A few more complex results

---

There are a lot of result that show that, under some (unfortunately strong) assumptions on  $\mathbf{X}$ , we can have consistency of  $\hat{\beta}_{\text{lasso},\lambda}$  or of its support when both  $n$  and  $p$  are large but the model is assumed to be truly sparse. See chapter 11 of of Hastie, Tibshirani & Wainwright (2015) for more on this.

There are also results that pretty much do not need any assumption. One of these is the one we will see in the next slide, that in particular does not assume that the "true model" is linear.

## An "assumption-free" result

---

We just assume that  $\mathbf{x}$  and  $y$  are two random variables, and that all features are bounded by some constant  $C$ . We **do not assume that  $\mathbf{x} \mapsto \mathbb{E}[y|\mathbf{x}]$  is linear.**

---

<sup>2</sup>various versions due to Greenshtein, Ritov, Juditsky, Nemirovski, and Wasserman, see more details on <https://normaldeviate.wordpress.com/2013/10/03/assumption-free-high-dimensional-inference/>

## An "assumption-free" result

---

We just assume that  $\mathbf{x}$  and  $y$  are two random variables, and that all features are bounded by some constant  $C$ . We **do not assume that  $\mathbf{x} \mapsto \mathbb{E}[y|\mathbf{x}]$  is linear**. For a parameter  $\beta$ , we define the predictive error

$$R(\beta) = \mathbb{E}_{\mathbf{x}, y} [(y - \mathbf{x}^T \beta)^2] .$$

---

<sup>2</sup>various versions due to Greenshtein, Ritov, Juditsky, Nemirovski, and Wasserman, see more details on <https://normaldeviate.wordpress.com/2013/10/03/assumption-free-high-dimensional-inference/>

# An "assumption-free" result

---

We just assume that  $\mathbf{x}$  and  $y$  are two random variables, and that all features are bounded by some constant  $C$ . We **do not assume that  $\mathbf{x} \mapsto \mathbb{E}[y|\mathbf{x}]$  is linear**. For a parameter  $\beta$ , we define the predictive error

$$R(\beta) = \mathbb{E}_{\mathbf{x}, y} [(y - \mathbf{x}^T \beta)^2].$$

Note that computing  $R(\beta)$  is impossible because we only have a finite data set. Let us consider

$$\beta^* \in \operatorname{argmin}_{\|\beta\|_1 \leq L} R(\beta).$$

The parameter  $\beta^*$  does not correspond to a "true model", but is the **best sparse linear predictor** that we could have computed with an infinite data set. **Theorem.<sup>2</sup> The lasso does almost as good as  $\beta^*$ :**

$$R(\hat{\beta}_{\text{lasso}, L}) \leq R(\beta^*) + \sqrt{\frac{8C^2L^4}{n} \log \left( \frac{2p^2}{\delta} \right)}.$$

---

<sup>2</sup>various versions due to Greenshtein, Ritov, Juditsky, Nemirovski, and Wasserman, see more details on <https://normaldeviate.wordpress.com/2013/10/03/assumption-free-high-dimensional-inference/>

## Let's recap: lasso theory

---

- **convex problem**, with very fast possible optimisation
- $\hat{\beta}_{\text{lasso},\lambda}$  is **sparse** (the larger the  $\lambda$ , the sparser)
- $\hat{\beta}_{\text{lasso},\lambda}$  will contain **at most  $\min(n, p)$  nonzero coefficients**