

Machine Learning

CS342

Lecture 7: Sparsity and Regularisation in Linear regression: The Lasso and Ridge Regression

Dr. Theo Damoulas

T.Damoulas@warwick.ac.uk

Office hours: Mon & Fri 10-11am @ CS 307



Model Sparsity and Interpretation

Sparsity in ML/Stats: *When few attributes (or few observations) contribute significantly to the resulting model. Many of them are “switched off”*

Why we might want to do that?

- **Statistical Interpretation:**
 - These attributes are more predictive than others
 - Discover significant correlations (e.g. income and education level)
 - Easier to “read”. Humans confuse correlation with causation as we tend to think causally. DTs very successful as easy to read.
- **Computational** reasons:
 - Resulting model is sparse so a reduced representation, less computation and less storage
 - Certain sparse models will become sparse even during training. So both training and prediction time improved (Big Data setting)
 - **Sparsity as a way to fight overfitting! Preference for simpler models**

Model Sparsity and Interpretation

Lets go back to linear regression with multiple attributes (multivariate)

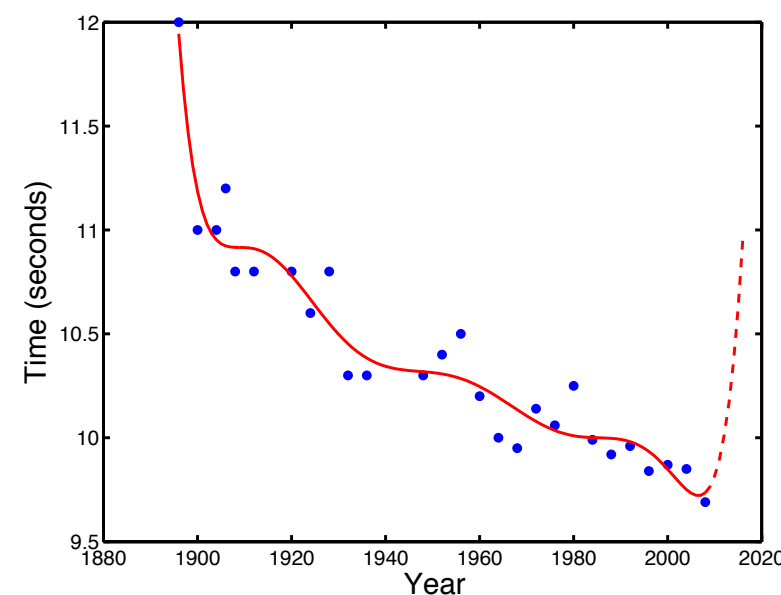
e.g. 2 attributes (2-D input space):

$$\hat{t}_n = \hat{w}_0 + \hat{w}_1 x_{n1} + \hat{w}_2 x_{n2} = \mathbf{x}_n \hat{\mathbf{w}}$$

e.g. 1 attribute but 4th order polynomial expansion (4-D input space):

$$\hat{t}_n = \hat{w}_0 + \hat{w}_1 x_{n1} + \hat{w}_2 x_{n1}^2 + \hat{w}_3 x_{n1}^3 + \hat{w}_4 x_{n1}^4 = \mathbf{x}_n \hat{\mathbf{w}}$$

$$\mathbf{X} = \begin{bmatrix} x_1^0 & x_1^1 & x_1^2 & \cdots & x_1^K \\ x_2^0 & x_2^1 & x_2^2 & \cdots & x_2^K \\ \vdots & \vdots & \vdots & \cdots & \vdots \\ x_N^0 & x_N^1 & x_N^2 & \cdots & x_N^K \end{bmatrix}$$





Model Sparsity and Interpretation

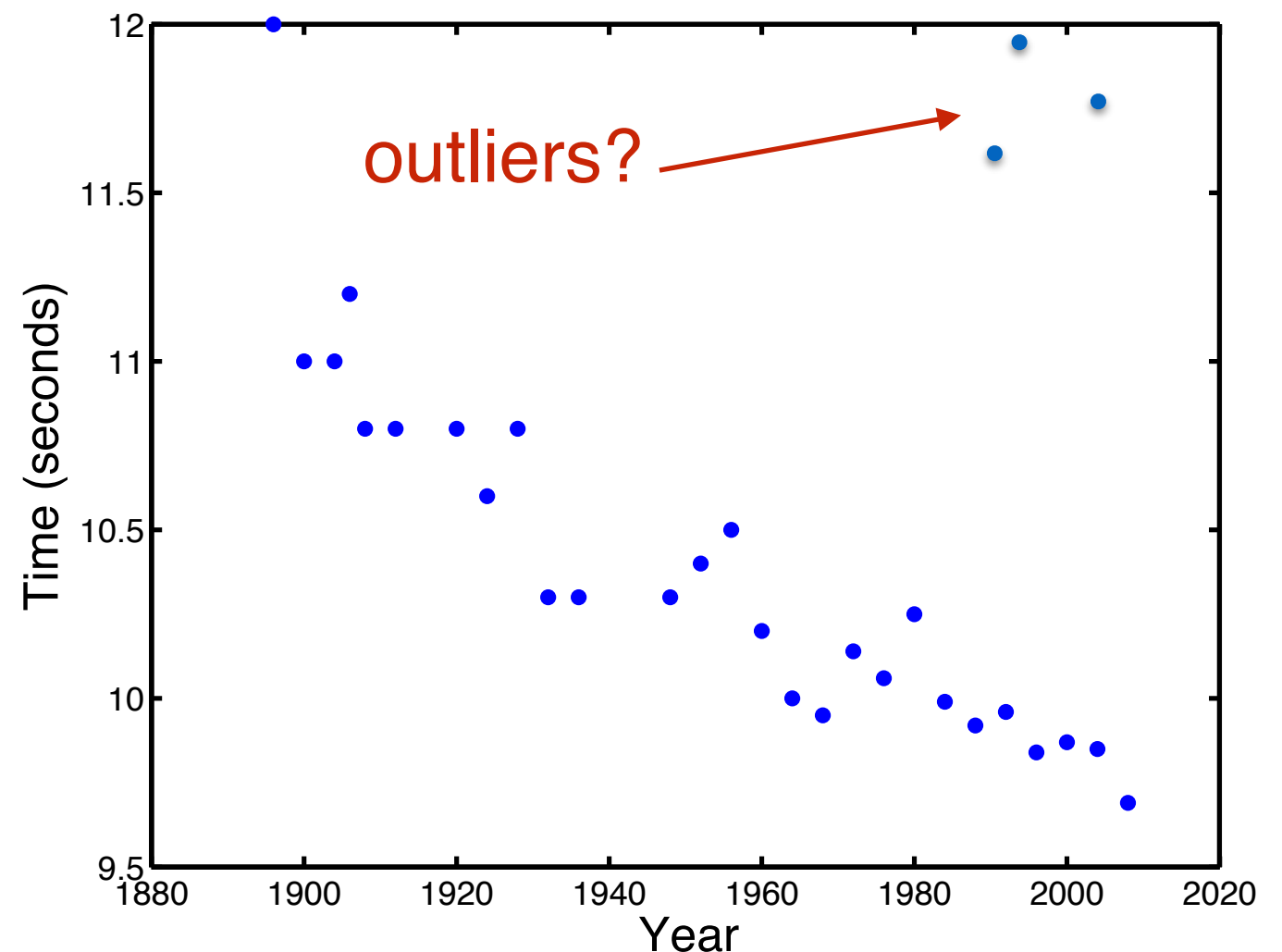
In both cases we did OLS (Minimise the squared error loss)

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

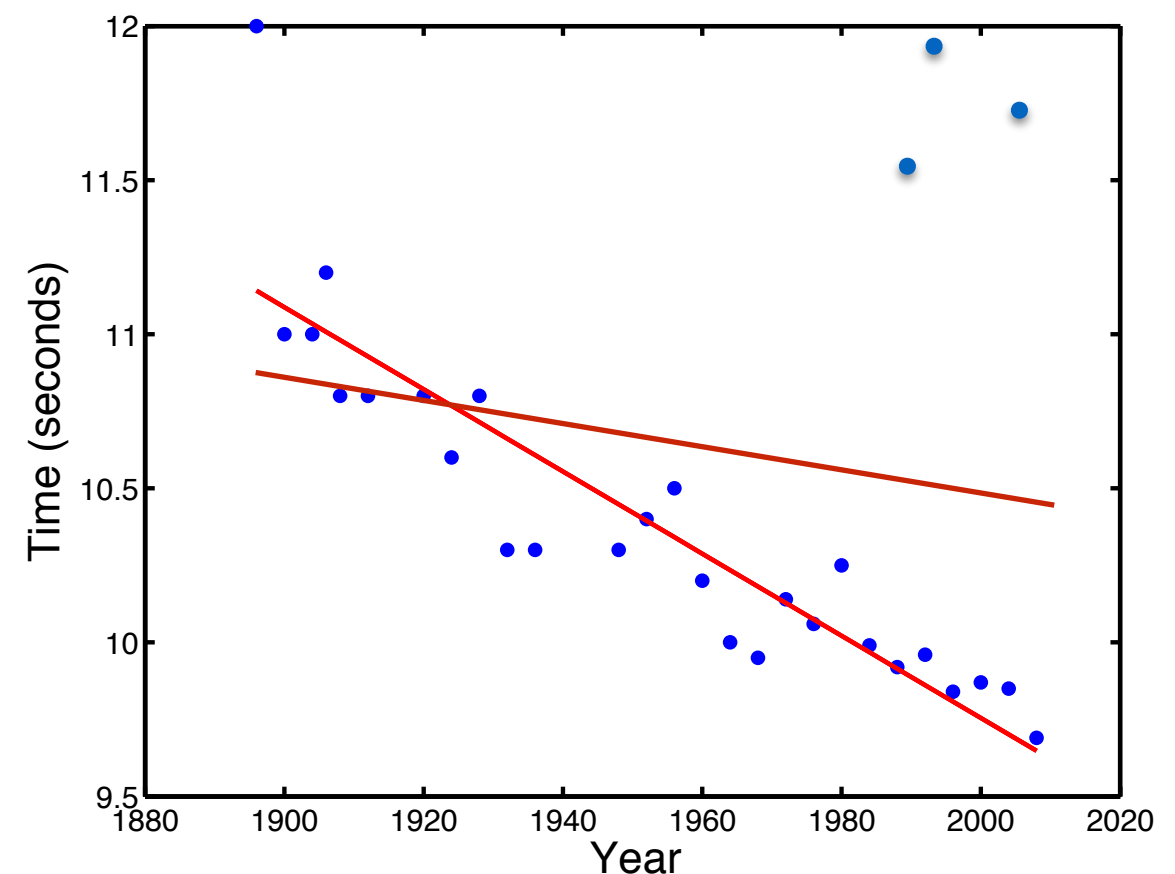
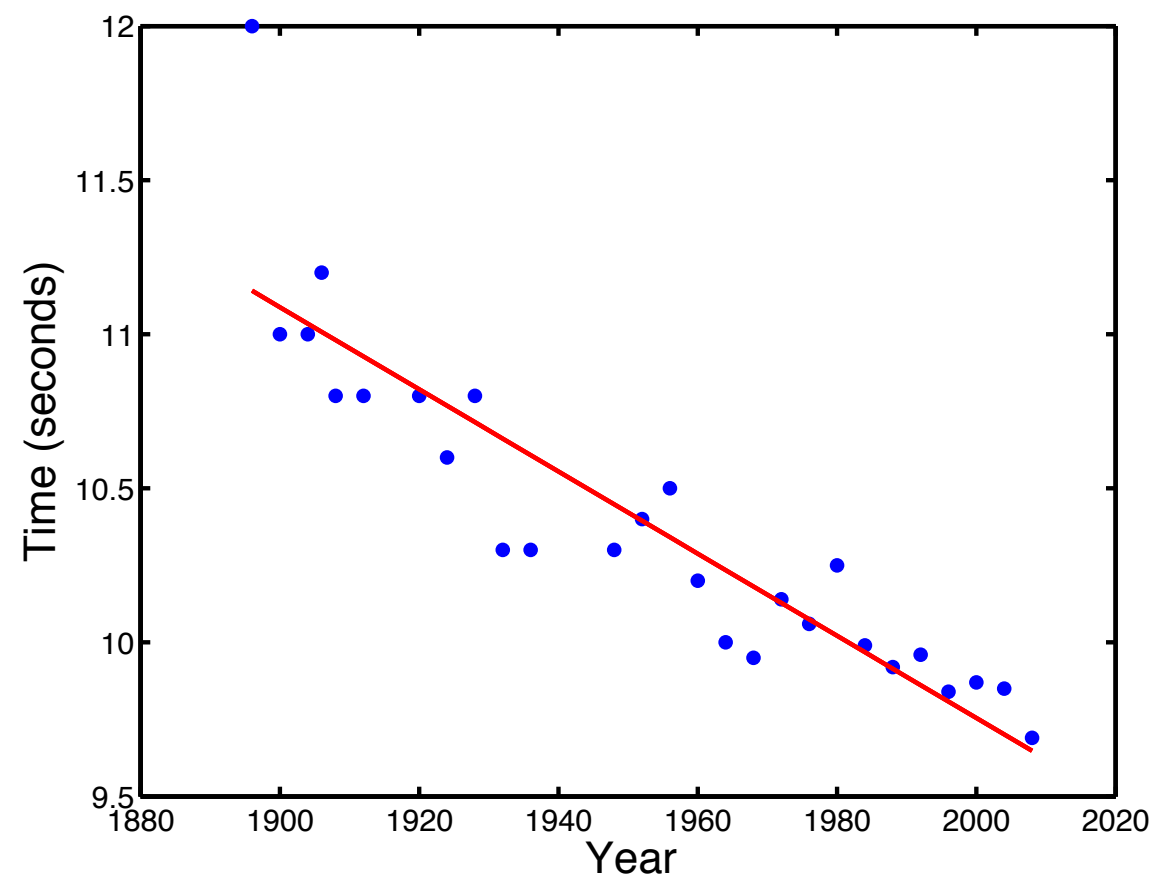
Our parameters are only “coupled” via the data
They can grow very large in magnitude - complex model

“Outliers”: *observations located far away from the rest of the data.*

They can dominate the OLS solution



Outliers strongly affect the OLS solution



And if we were fitting here a polynomial expansion with OLS
we would get a complex model

So what can we do to “penalise very complex solutions”?

Regularisation: Coupling our parameters and Penalising

We need to encode our preference for simpler solutions!
In regularisation we do that by coupling the parameter magnitudes

We impose a “competition” between them and constraint their magnitudes.

Rogers & Girolami Ch. 1, p. 33-34

One way: keep this value low $\sum_d w_d^2 = \mathbf{w}^T \mathbf{w}$

Update our Loss function to include it $\mathcal{L}' = \mathcal{L} + \lambda \mathbf{w}^T \mathbf{w}$

Why this value? Why adding it to Loss function?

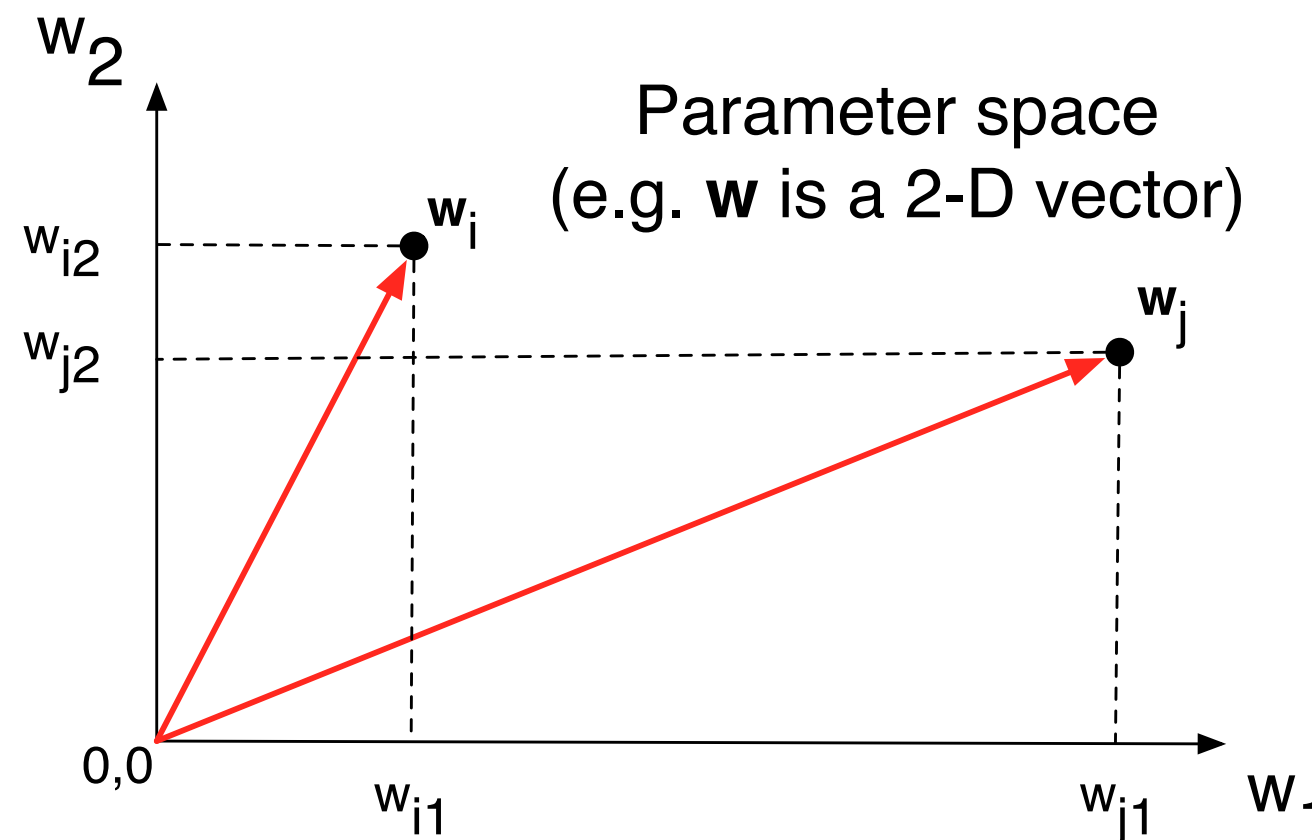
What is this regulariser we have added to the Loss?

$$\sum_d w_d^2 = \mathbf{w}^T \mathbf{w}$$

The square of the Euclidean distance of \mathbf{w} vector from centre of axis

Wait... isn't that the same as the square of the L_2 norm?

Euclidean (L_2) Distance of \mathbf{w} from (0,0)
= L_2 norm of \mathbf{w}





Regularisation as adding Constraints

Original OLS Loss :
(vector-matrix format)

$$\mathcal{L} = \frac{1}{N} (\mathbf{t} - \mathbf{X}\mathbf{w})^T (\mathbf{t} - \mathbf{X}\mathbf{w})$$

I still want to minimise this loss **but subject to constraints**

constraint: “Keep parameters small by coupling their magnitudes”
or... The (square of the) L2 norm of \mathbf{w} stays within a limit

$$\text{Minimise } \mathcal{L} \quad \text{s.t.} \quad \sum_d w_d^2 = \mathbf{w}^T \mathbf{w} \leq t$$

Equivalent to minimising $\mathcal{L}' = \mathcal{L} + \lambda \mathbf{w}^T \mathbf{w}$



Penalised Least Squares: Ridge Regression

Known as: Regularised LS - Penalised LS - Ridge Regression

Repeat derivation for new solution (PLS)

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X} + N \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{t}$$

But what does the added constraint say in geometrical terms?

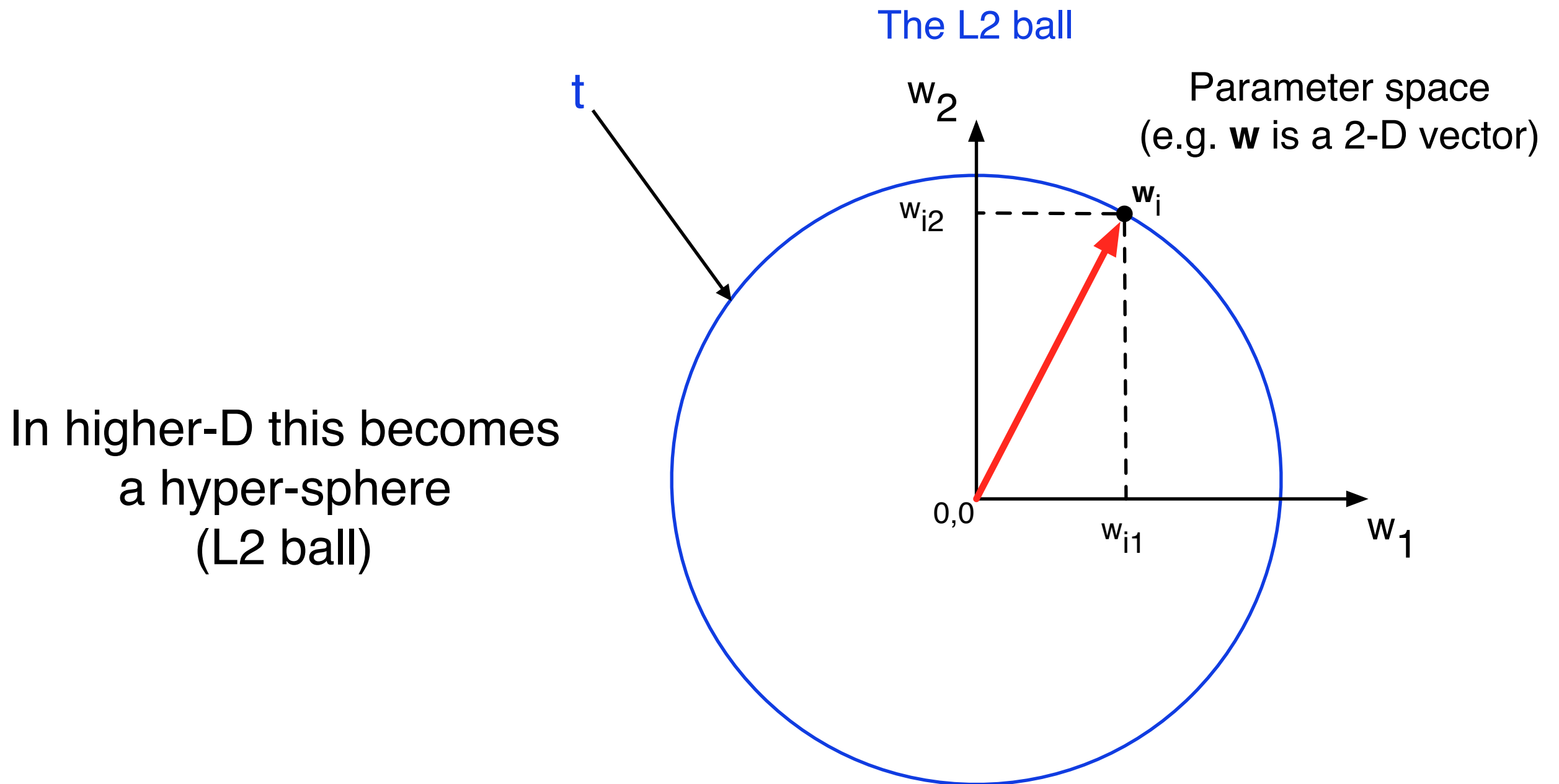
$$L_2^2(\mathbf{w}) = \sum_d w_d^2 = \mathbf{w}^T \mathbf{w} \leq t$$

Square of Euclidean distance (L_2) of \mathbf{w} from (0,0) within a limit t

Lets forget about the square for now..



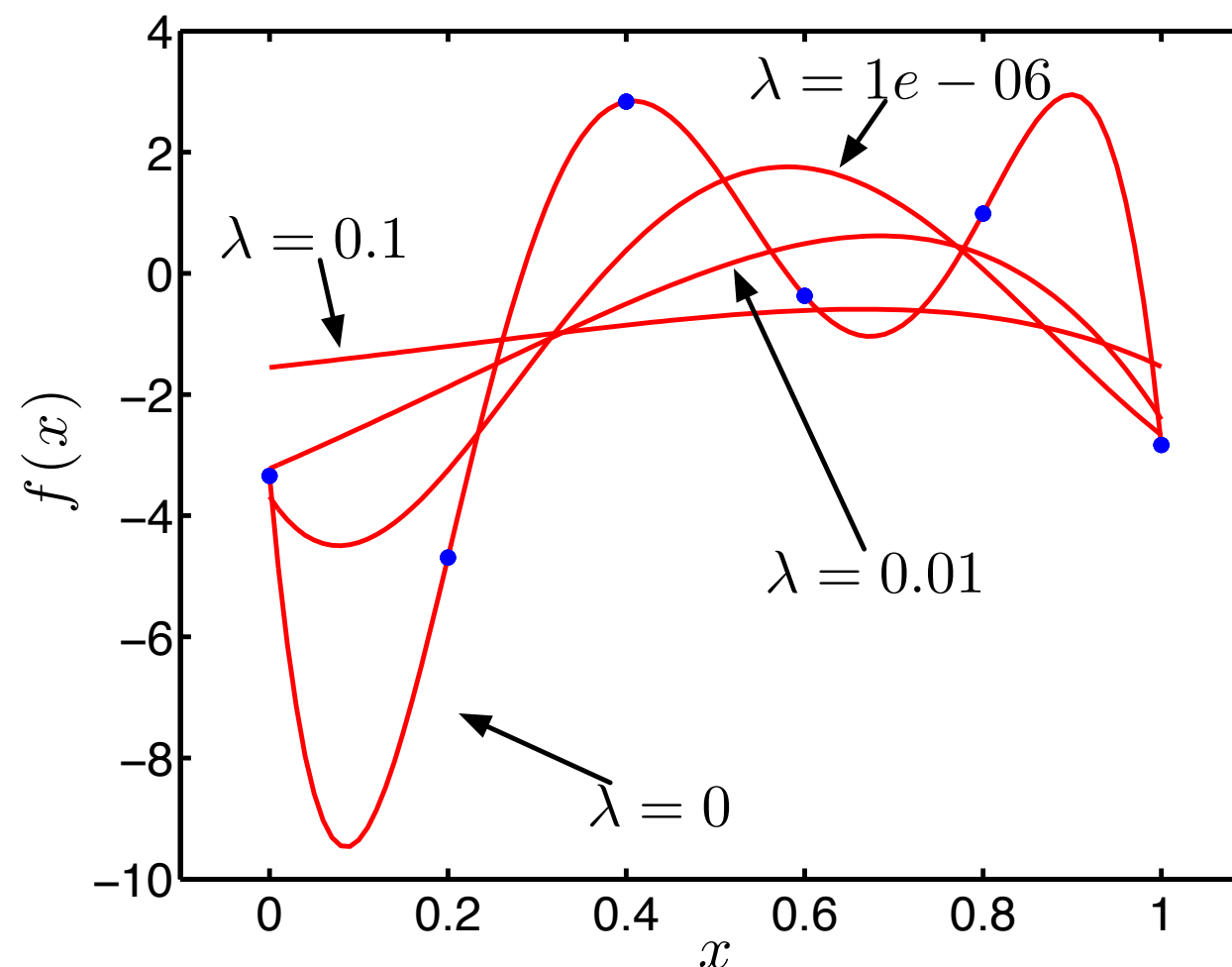
Geometry: the L2 ball and hyper-spheres



Penalised Least Squares - Ridge Regression

6 datapoints, 5-order polynomial, varying the strength of regularisation

You can use CV to choose lambda, the strength of penalisation (complexity)





Summary so far

In Statistics: parameters = Regression coefficients

OLS

- Linear regression with OLS can easily overfit complex models
- Linear regression with OLS can be mislead by outliers
- The magnitudes of the parameters \mathbf{w} in OLS are not constrained
- Numerical instabilities (high variance) when correlated attributes

PLS / Ridge regression

- We “couple” the parameter magnitudes to constrain them
- We do that by adding a regulariser to the minimisation problem
- In PLS that regulariser is the square of the L_2 norm $L_2^2(\mathbf{w})$
- This is also known as “Ridge Regression”
- The solution becomes: $\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X} + N\lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{t}$
- Lambda controls the strength of regularisation (the volume of the ball)

If you change regulariser you get other algorithm(s)

The Lasso



<http://statweb.stanford.edu/~tibs/lasso.html>

Use another “stronger” regulariser instead of (a function of) the L_2
Follow an algorithmic procedure to shrink some \mathbf{w} 's to 0

The Lasso

L_2 norm does not really guarantee “sparsity”

It will constraint the parameter magnitudes but will not necessarily set some to 0

L_1 norm is a stronger penaliser and will force parameters (of irrelevant/uncorrelated attributes) to 0

Geometry: L_2 ball (hyper-sphere). What is the L_1 norm?

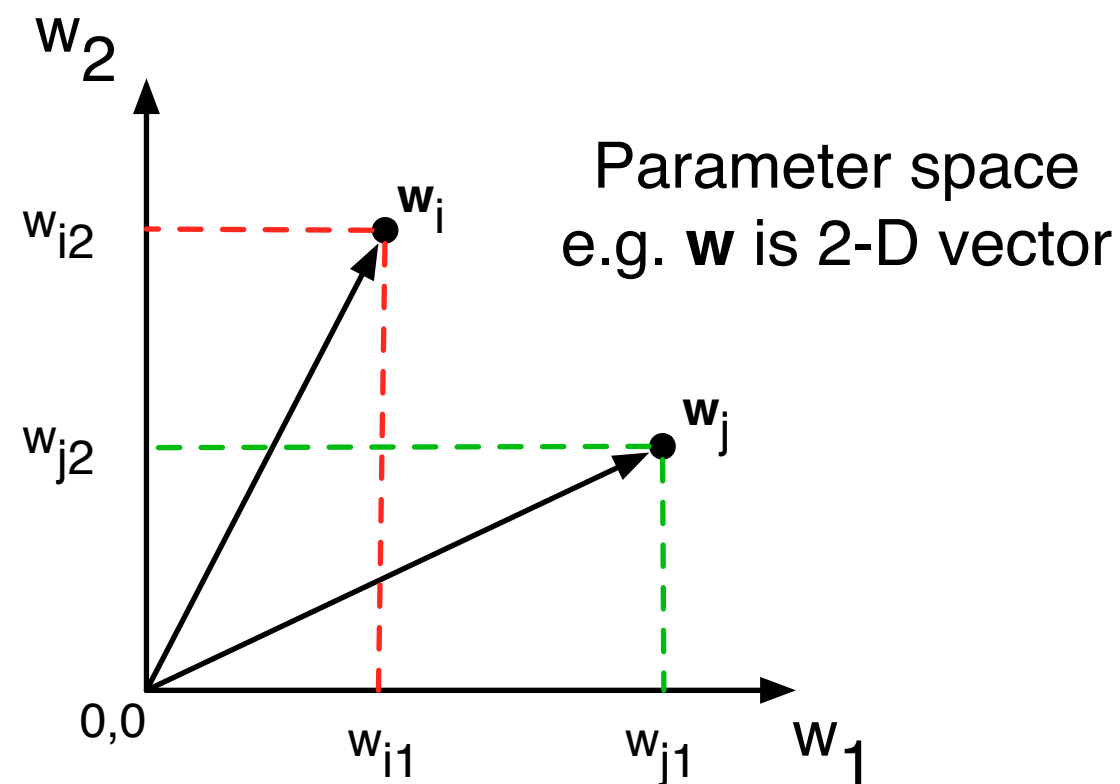
L_p definition

$$L_p(\mathbf{w}) = \left(\sum_{d=1}^D |w_{nd}|^p \right)^{\frac{1}{p}}$$

L₁ norm

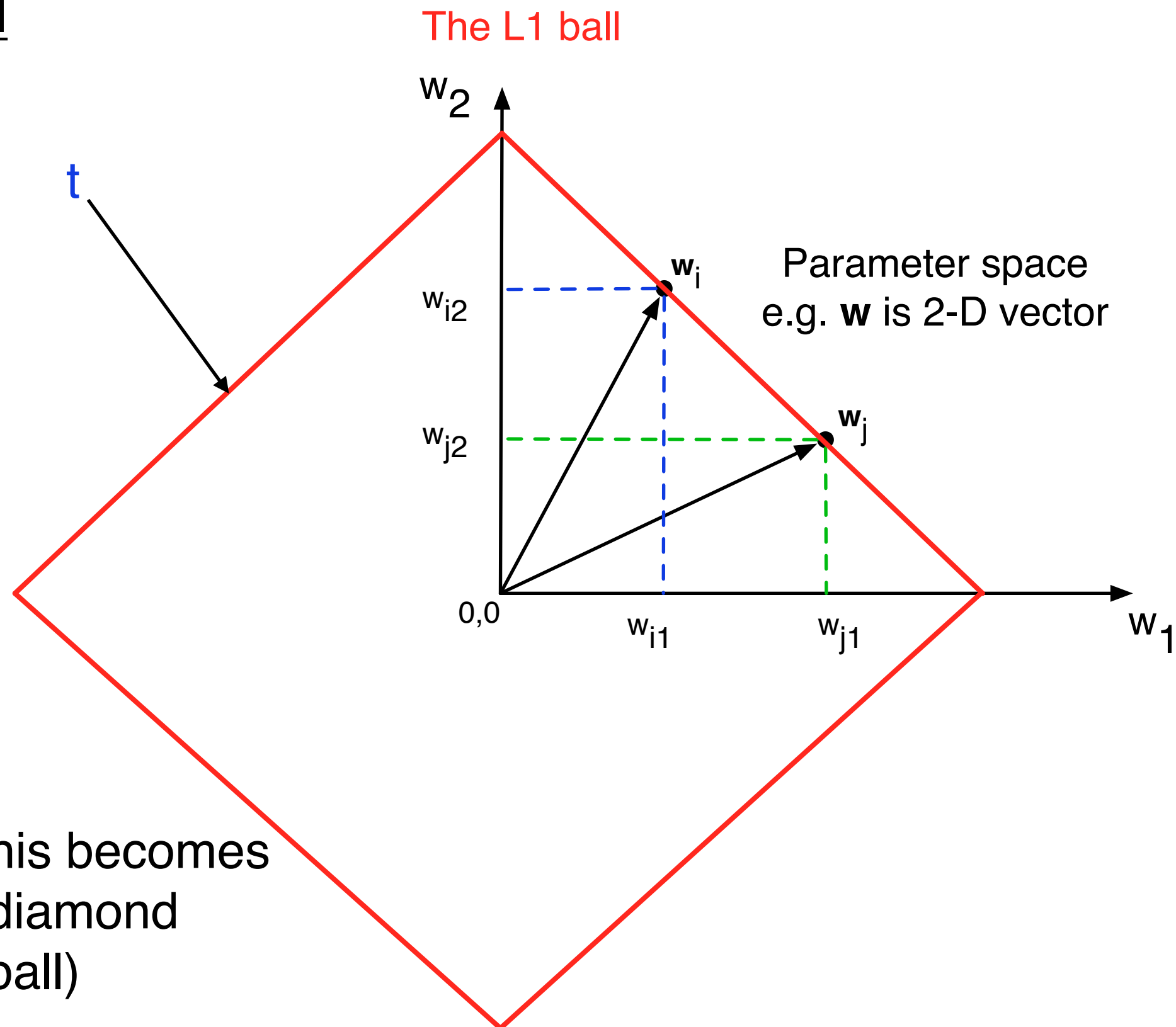
$$L_1(\mathbf{w}) = \sum_{d=1}^D |w_d|$$

Manhattan (L1) Distance of \mathbf{w} from (0,0)
= L1 norm of \mathbf{w}



So if I keep that distance constant and draw all the \mathbf{w} 's
What shape will I get?

The L_1 ball



Some of the L_p balls

Volume of L_p balls in 3-D for $L_p \leq 1$



$p = \infty$



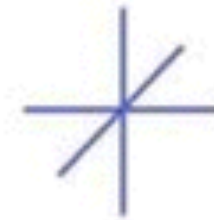
$p = 2$



$p = 1$



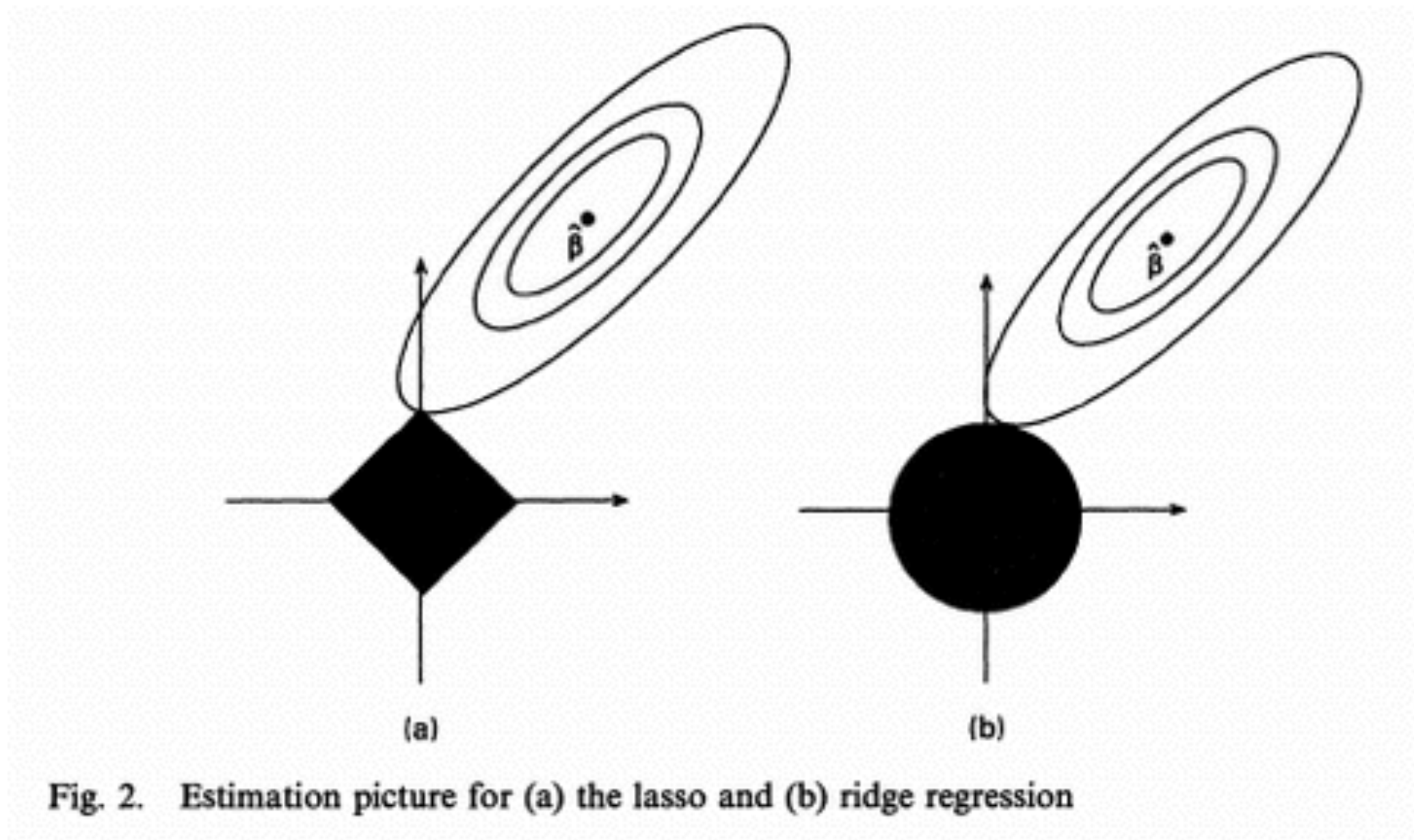
$0 < p < 1$



$p = 0$

Regularisation with L_2 is also known as Tikhonov regularisation

The Lasso vs Ridge Regression (PLS)



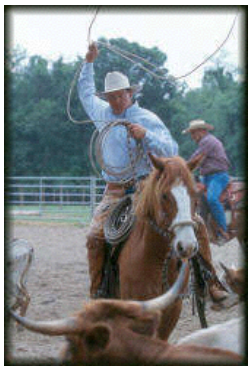
The Lasso will recover *sparse* solutions where some parameters will be 0. Hence the corresponding attribute will be ignored



The Lasso

Original OLS Loss :
(vector-matrix format)

$$\mathcal{L} = \frac{1}{N} (\mathbf{t} - \mathbf{X}\mathbf{w})^T (\mathbf{t} - \mathbf{X}\mathbf{w})$$



Minimise \mathcal{L} s.t. $\sum_d |w_d| \leq t$

Equivalent to minimising $\mathcal{L}' = \mathcal{L} + \lambda \sum_d |w_d|$

Quadratic programme (optimisation)

Use **cross-validation** to pick strength of regularisation lambda



Summary of The Lasso

Do not use OLS like most social scientists do!

PLS / Ridge regression

`sklearn.linear_model.Ridge`

- We “couple” the parameter magnitudes to constrain them
- We constraint parameters by regularising with squared **L₂ norm**
- Lambda controls the strength of regularisation (the volume of the ball)

The Lasso

`sklearn.linear_model.Lasso`

- We “couple” the parameter magnitudes to constrain them
- We constraint parameters by regularising with the **L₁ norm**
- Sparse solutions with some parameters at 0
- Great for Interpretation
- Will under-fit if our problem is not really sparse (use PLS instead)
- Will outperform PLS when many attributes are irrelevant

Other variants (Elastic Net) with mixed norms!