

CS5014 Machine Learning

Lecture 10 Regularisation

Lei Fang

School of Computer Science, University of St Andrews

Spring 2021



University
of
St Andrews

Some responses

What one hot encoding actually does ?

- regression ?
- classification ?

Why Newton's step use Hessian's inverse \mathbf{H}^{-1} ?

$$\boldsymbol{\theta}_{t+1} \leftarrow \boldsymbol{\theta}_t - \mathbf{H}_t^{-1} \mathbf{g}_t$$

- long story: it optimise local approximated quadratic function
- short story: inversely related to curvature (H measures curvature in higher dimensions)
 - larger curvature (pointy and curvy) \Rightarrow shorter steps;
 - smaller curvature (flatter)

Today's topic

Nonlinear model

- fixed basis models

Regularisation

- l_2 penalty: ridge regression
- l_1 penalty: lasso regression

Towards nonlinear models: regression

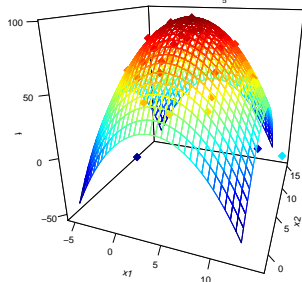
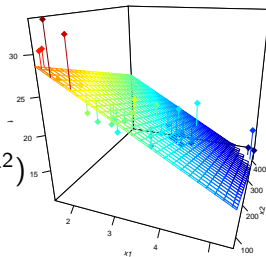
For linear regression:

$$P(y|\mathbf{x}, \boldsymbol{\theta}) = N(\underbrace{f(\mathbf{x}; \boldsymbol{\theta})}_{\boldsymbol{\theta}^T \mathbf{x}: \text{linear}}, \sigma^2) \text{ or } y = \underbrace{f(\mathbf{x}; \boldsymbol{\theta})}_{\boldsymbol{\theta}^T \mathbf{x}} + \epsilon, \quad \epsilon^{(i)} \sim N(0, \sigma^2)$$

The regression function f is assumed linear

$$f(\mathbf{x}; \boldsymbol{\theta}) = \boldsymbol{\theta}^T \mathbf{x}$$

- *i.e.* fitting lines/hyperplanes
- in real life, a lot of relationships are not linear
- and we do not know what $f(\mathbf{x})$ should look like !



Towards nonlinear models: classification

For logistic regression:

$$P(y|\mathbf{x}, \boldsymbol{\theta}) = \text{Ber}(\underbrace{\sigma(f(\mathbf{x}; \boldsymbol{\theta}))}_{\boldsymbol{\theta}^T \mathbf{x}: \text{linear}}) = \sigma^y (1 - \sigma)^{(1-y)}$$

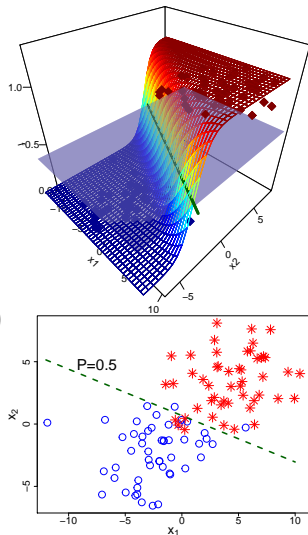
To predict the label y for any input \mathbf{x} :

$$y = 1 \text{ if } P(y = 1|\mathbf{x}, \boldsymbol{\theta}) > 0.5 \quad y = 0 \text{ if otherwise}$$

Note that the **decision boundary** is linear (hyperplane or line)

$$P(y = 1|\mathbf{x}, \boldsymbol{\theta}) = 0.5 \Rightarrow \boldsymbol{\theta}^T \mathbf{x} = 0$$

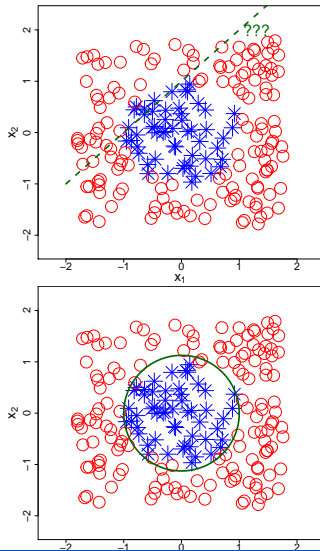
- *i.e.* separating data by lines/hyperplanes
- in reality, we do know what f should be; plane or a more general surface



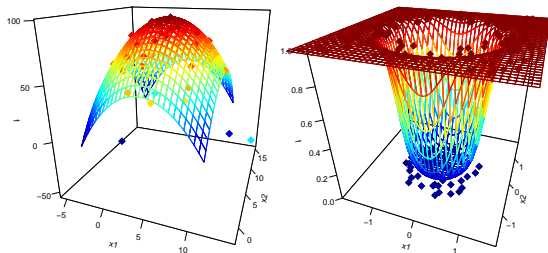
Nonlinear classification data

What if your data looks like this ?

- no linear decision boundary or $f(\mathbf{x}; \boldsymbol{\theta}) = \boldsymbol{\theta}^T \mathbf{x}$ seems making much sense
- but a non-linear boundary makes more sense
 - the classification rule is actually $\|\mathbf{x}\|_2^2 = x_1^2 + x_2^2 \leq 1$
 - distance to $\mathbf{0}$ is less than 1
 - the boundary is a circle
 - I know it because I generated the data



Nonlinear model: polynomial model



Both models are actually 2nd order polynomial:

$$f(\mathbf{x}; \beta) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2 + \beta_4 x_1^2 + \beta_5 x_2^2 = \sum_{k=0}^5 \beta_k \phi_k(\mathbf{x}) = \beta^T \phi(\mathbf{x})$$

- where

$$\phi(\mathbf{x}) = [\underbrace{1}_{\phi_0(\mathbf{x})}, \underbrace{x_1}_{\phi_1(\mathbf{x})}, \underbrace{x_2}_{\phi_2(\mathbf{x})}, \underbrace{x_1 x_2}_{\phi_3(\mathbf{x})}, \underbrace{x_1^2}_{\phi_4(\mathbf{x})}, \underbrace{x_2^2}_{\phi_5(\mathbf{x})}]^T$$

- it expands $\mathbf{x} = [1, x_1, x_2]^T$ to a larger vector

Nonlinear response from linear model

Note that you get a free nonlinear model by transforming the input \mathbf{X}

$$\mathbf{X} = \begin{bmatrix} 1 & x_1^{(1)} & x_2^{(1)} \\ 1 & x_1^{(2)} & x_2^{(2)} \\ \vdots & \vdots & \vdots \\ 1 & x_1^{(m)} & x_2^{(m)} \end{bmatrix} \Rightarrow \Phi = \begin{bmatrix} \phi(\mathbf{x}^{(1)}) \\ \phi(\mathbf{x}^{(2)}) \\ \vdots \\ \phi(\mathbf{x}^{(m)}) \end{bmatrix} = \begin{bmatrix} 1 & x_1^{(1)} & x_2^{(1)} & x_1^{(1)}x_2^{(1)} & (x_1^{(1)})^2 & (x_2^{(1)})^2 \\ 1 & x_1^{(2)} & x_2^{(2)} & x_1^{(2)}x_2^{(2)} & (x_1^{(2)})^2 & (x_2^{(2)})^2 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_1^{(m)} & x_2^{(m)} & x_1^{(m)}x_2^{(m)} & (x_1^{(m)})^2 & (x_2^{(m)})^2 \end{bmatrix}$$

- remember superscript (i) index data samples; and subscript index features
- Φ is a $m \times 6$ matrix

The expanded model for regression is

$$\mathbf{y} = \Phi\beta + \epsilon$$

- still a linear model w.r.t ϕ , the expanded new features
- all existing results apply: gradient descent, normal equation (replace \mathbf{X} with Φ)

$$\beta_{ls} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y}$$

Polynomial basis function expansion

Consider predictors with $n = 2$ input features,

$$\mathbf{x} = [1, x_1, x_2]^T \text{ (1 is dummy variable)}$$

Second order polynomial, ϕ expands \mathbf{x} to

$$\phi(\mathbf{x}) = [1, x_1, x_2, x_1x_2, x_1^2, x_2^2]^T;$$

- what if the input size is n rather than 2? $(1 + n + \binom{n}{2} + n) \in O(n^2)$

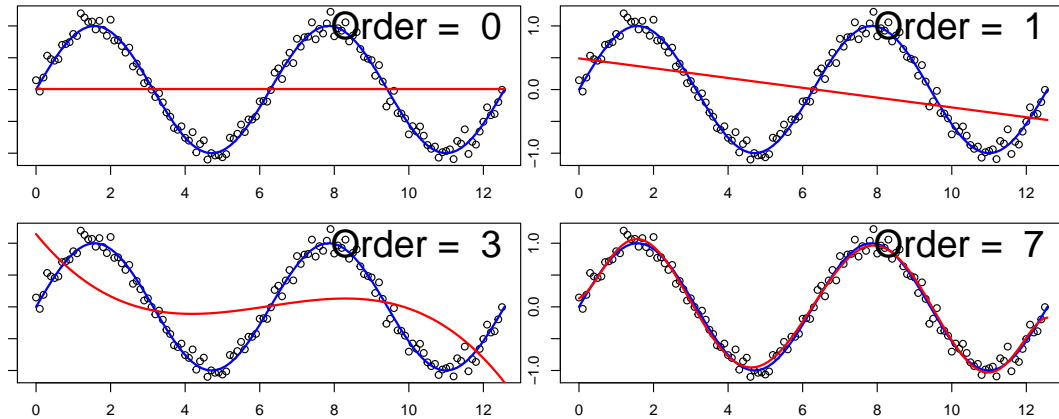
Third order polynomial, it becomes

$$\phi(\mathbf{x}) = [1, x_1, x_2, x_1x_2, x_1^2, x_2^2, \underbrace{x_1x_2^2, x_1^2x_2, x_1^3, x_2^3}_{\text{3-rd order terms}}]^T$$

Higher order leads to larger number of basis; it may not necessary be a good thing though

Example: polynomial fitting

True regression function: $f = \sin(x)$; and $y^{(i)} = \sin(x^{(i)}) + \epsilon^{(i)}$, $\epsilon \sim N(0, 0.2^2)$



Other basis function ϕ

$\phi(\mathbf{x})$ can take other forms

- each $\phi_k(\mathbf{x})$ is a $R^n \rightarrow R$ transformation;
 - so ϕ is a $R^n \rightarrow R^p$ transformation: previous example: $n = 3$, $p = 6$
 - obviously, if $\phi(\mathbf{x}) = \mathbf{I}\mathbf{x}$, we recover ordinary linear regression
 - ▶ it is a specific case of basis expansion model
- 2nd order polynomial: $\phi_k(\mathbf{x}) = x_j^2$, $x_j x_{j'}$, x_j , or 1; for $k = 1 \dots p$
- $\phi_k(x) = \log(x)$, \sqrt{x} can take other (literally any) nonlinear forms

Fixed basis models

Fixed basis model assumes

$$f(\mathbf{x}; \boldsymbol{\beta}) = \sum_{k=0}^{p-1} \beta_k \phi_k(\mathbf{x}) = \boldsymbol{\beta}^T \boldsymbol{\phi}(\mathbf{x})$$

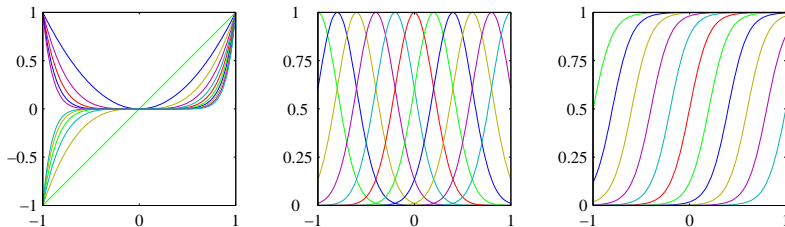
- “fixed” because the basis functions need to be manually specified
 - the contrary is adaptive basis (the basis are learnt)
- still linear w.r.t each new feature input ϕ_k
 - all linear learning results apply to $\boldsymbol{\beta}$
- no longer linear in the original input \mathbf{x}

Two popular basis in ML literature

In traditional Machine Learning literature, radial basis function (RBF) and sigmoid (or “tanh”) are popular

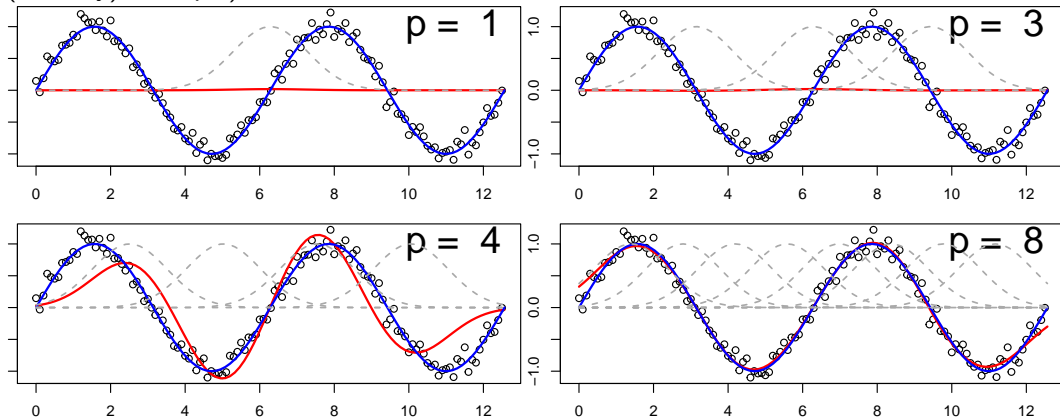
$$\underbrace{\phi_k(\mathbf{x}) = \exp\left(-\frac{\|\mathbf{x} - \boldsymbol{\mu}_k\|_2^2}{2s^2}\right)}_{\text{radial basis}}, \quad \underbrace{\phi_k(x) = \frac{1}{1 + \exp\left(-\frac{x - \mu_k}{s}\right)}}_{\text{sigmoid basis}}$$

- they are location based: depending on μ ; while polynomial basis are not (lower left)!
- very general basis and can fit various different models



Example: radial basis function (RBF) with varying number of basis

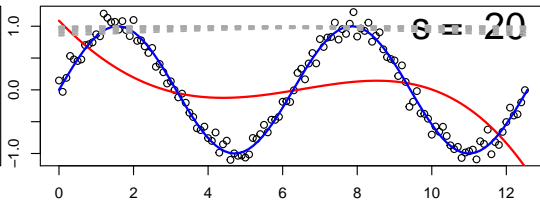
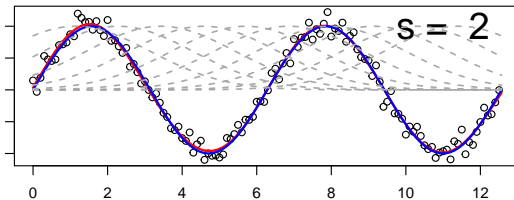
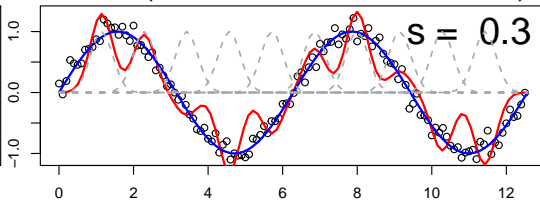
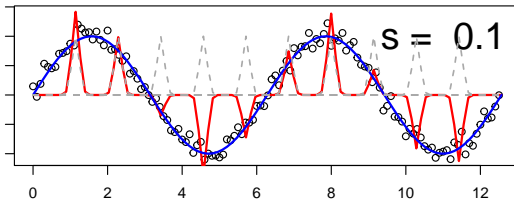
True regression function: $\mathbf{f} = \sin(\mathbf{x})$; and $y^{(i)} = \sin(x^{(i)}) + \epsilon^{(i)}$, $\epsilon \sim N(0, 0.2^2)$
More RBF basis (dashed gray lines) fits better; location μ_k matter ($p = 3$ is an (unlucky) example);



Example: radial basis function (RBF) with varying scale s

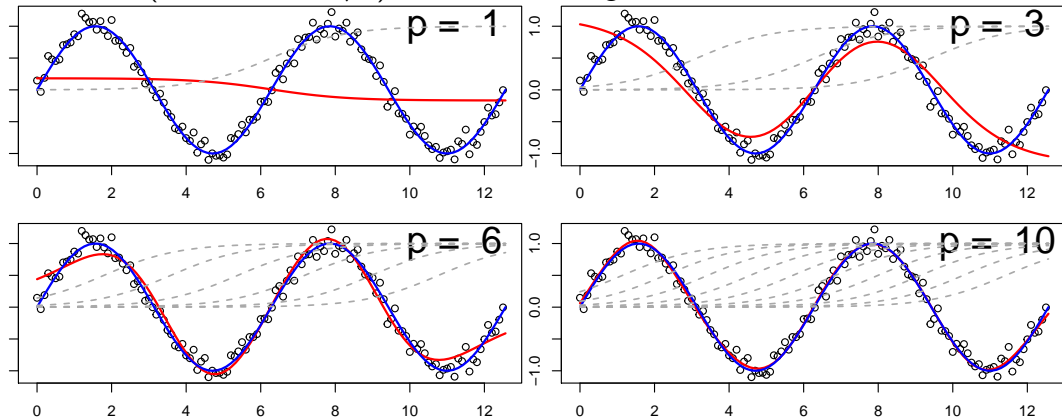
Smaller scale $s \Rightarrow$ wiggly predictions;

Large scale $s \Rightarrow$ flatter predictions, not good either ! ($p = 10$ basis for all four cases)



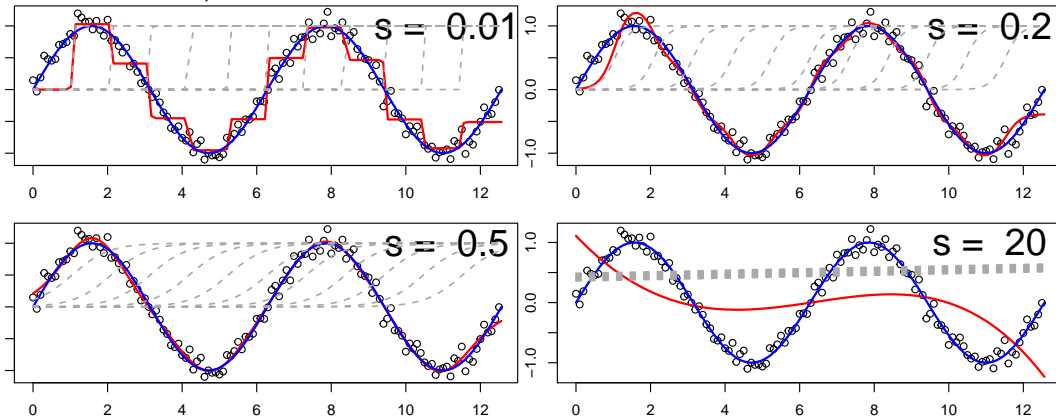
Example: sigmoid basis with varying number of basis

More basis (more locations μ_k) fits data better again



Example: sigmoid basis with varying scale s

Smaller scale $s \Rightarrow$ step functions; Large scale $s \Rightarrow$ smoother predictions ($p = 11$ basis for all four cases)



Interpretation of fixed basis models

In abstract vector space, functions are vectors, say $\phi_k(\mathbf{x})$

$$f(\mathbf{x}; \beta) = \sum_{k=0}^{p-1} \beta_k \phi_k(\mathbf{x})$$

- f : a linear combination of vectors $\{\phi_0, \phi_1, \dots, \phi_p\}$
- e.g. $f(x) = \beta_0 + \beta_1 x$ is a linear combination of two functions (vectors): $\phi_0(x) = 1$ and $\phi_1(x) = x$

$$f(x) = \beta_0 \phi_0(x) + \beta_1 \phi_1(x)$$

- linear regression: $\hat{f} \in \text{span}(\{1, x\})$ with $\hat{\beta}_0, \hat{\beta}_1$
 - projection of f to the subspace
- larger functional space $\text{span}(\{\phi_k(x)\}_1^p)$ (p increases) \Rightarrow better fit
 - until f exactly lives in the span (which can be bad!)
 - regardless of ϕ_k 's shapes as long as nonlinear and “different” (linear independent)

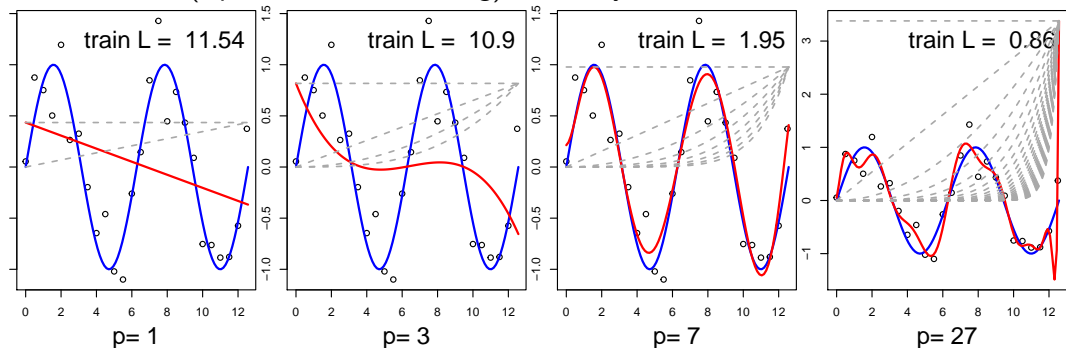
Flexibility resolves underfitting but introduce overfitting

Training loss:

$$L(\beta) = \sum_{i=1}^m (y^{(i)} - \beta^T \phi(\mathbf{x}^{(i)}))^2$$

More basis functions (p increases) \Rightarrow flexible model (larger functional subspace)

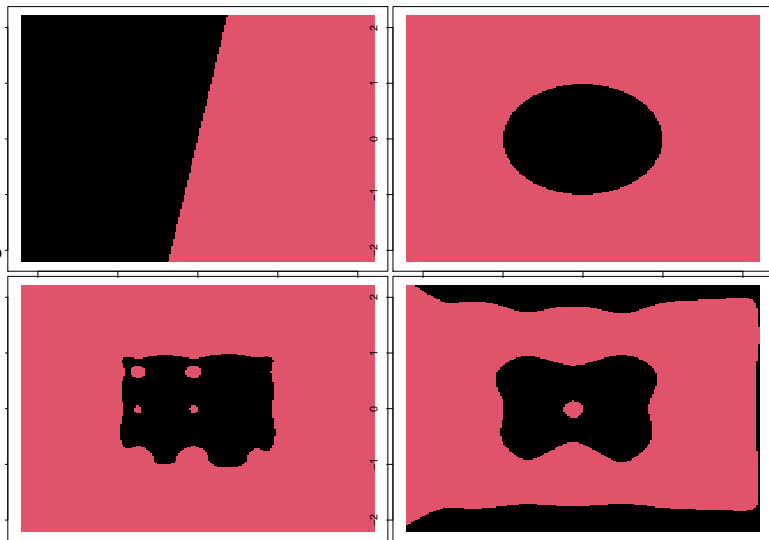
- MLE loss (squared error in training) uniformly decreases



Classification example: polynomial basis

The training NLL drops

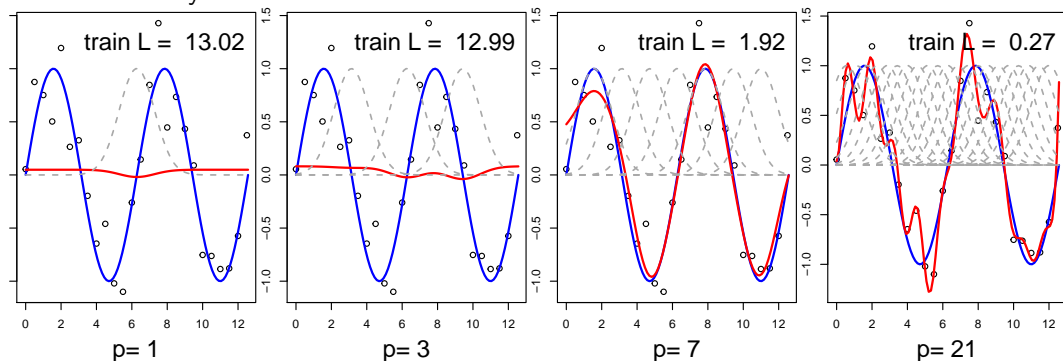
- possible to achieve 100% training accuracy
- regardless of dataset
- just keep expanding



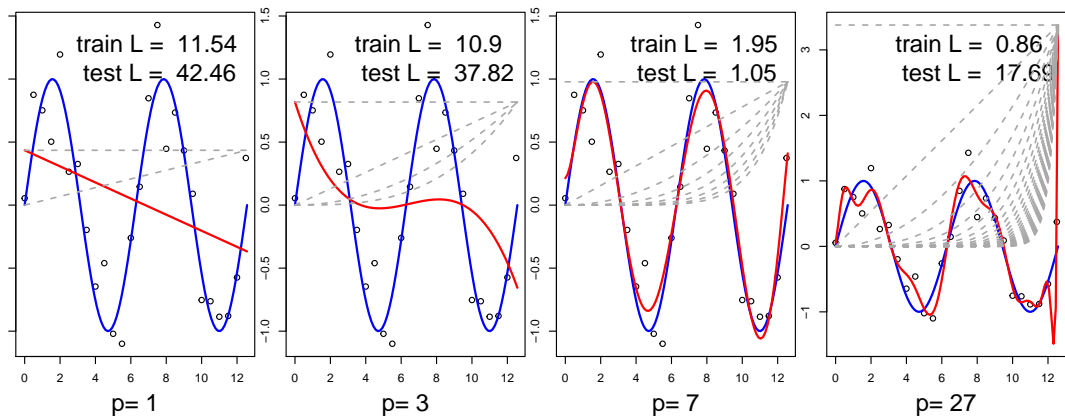
Flexibility resolves underfitting but introduce overfitting (RBF example)

The problem applies to all basis (no matter what shape)

- as long as you expands the basis \Rightarrow smaller training error
- MLE always favour flexible models



Testing results reveals severe overfitting



Regularisation: ridge regression

For linear regression, we add a l_2 norm penalty to the loss

$$L_{\text{ridge}}(\beta) = \underbrace{\sum_{i=1}^m (y^{(i)} - \beta^T \phi(\mathbf{x}^{(i)}))^2}_{\text{previous NLL loss}} + \lambda \underbrace{\beta^T \beta}_{L_2 \text{ penalty}}$$

- as NLL (squared error loss) term always favour flexible models
- we add a penalty term: $\beta^T \beta = \|\beta\|_2^2$; its l_2 norm
- large $\beta \Rightarrow$ higher penalty
 - “regularise” the value of β : large value of β is discouraged
- large $\lambda \Rightarrow$ large penalty: controls the penalty

Learning of ridge regression

Ridge regression has closed form solution (no surprise: still a quadratic function of β):

$$\beta_{\text{ridge}} = \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^m (y^{(i)} - \beta^T \phi(\mathbf{x}^{(i)}))^2 + \lambda \beta^T \beta = \underset{\beta}{\operatorname{argmin}} \underbrace{(\mathbf{y} - \Phi \beta)^T (\mathbf{y} - \Phi \beta) + \lambda \beta^T \beta}_{L_{\text{ridge}}(\beta)}$$

Take the derivative

$$\nabla_{\beta} L_{\text{ridge}} = -2(\mathbf{y} - \Phi \beta)^T \Phi + 2\lambda \beta^T \quad (\text{note } \nabla \beta^T \beta = 2\beta^T)$$

And set it to zero

$$\begin{aligned} -2(\mathbf{y} - \Phi \beta)^T \Phi + 2\lambda \beta^T &= \mathbf{0} \Rightarrow \Phi^T \Phi \beta + \lambda \beta = \Phi^T \mathbf{y} \\ \Rightarrow (\Phi^T \Phi + \lambda \mathbf{I}) \beta &= \Phi^T \mathbf{y} \quad (\text{note } \lambda \beta = \lambda \mathbf{I}_{p \times p} \beta) \\ \Rightarrow \beta_{\text{ridge}} &= (\Phi^T \Phi + \lambda \mathbf{I})^{-1} \Phi \mathbf{y} \end{aligned}$$

Learning for logistic regression (with l_2 penalty)

Similarly, we can derive the results for logistic regression (l_2 penalty)

$$\beta_{\text{ridge}} = \underset{\beta}{\operatorname{argmin}} \underbrace{- \left(\sum_{i=1}^m y^{(i)} \log \sigma^{(i)} + (1 - y^{(i)}) \log(1 - \sigma^{(i)}) \right)}_{\text{Negative log likelihood}} + \underbrace{\lambda \beta^T \beta}_{l_2 \text{ penalty}}$$

The gradient is

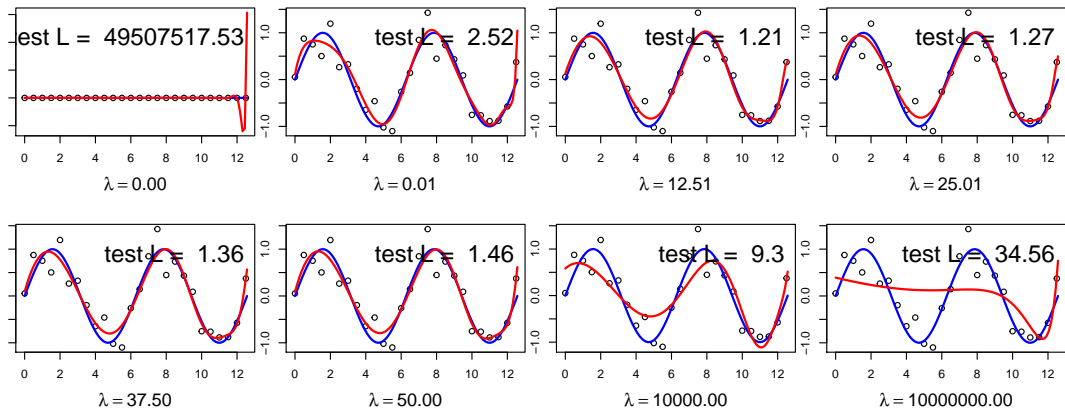
$$\nabla_{\beta} L_{\text{ridge}} = \underbrace{-(\mathbf{y} - \sigma(\Phi\beta))^T \Phi}_{\text{same as before but negated}} + \underbrace{2\lambda\beta^T}_{\nabla \text{penalty}} \quad (\text{I am abusing the notation for } \sigma \text{ as before})$$

The Hessian is

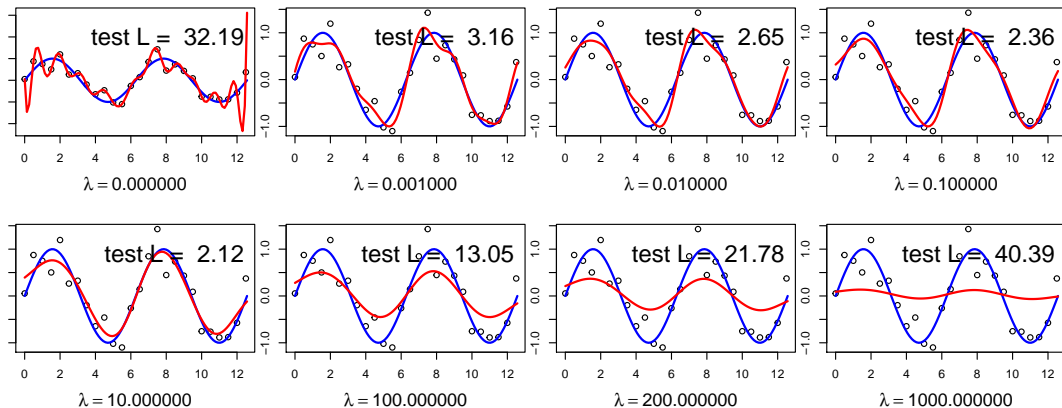
$$H_{\text{ridge}} = -\mathbf{X}^T \mathbf{D} \mathbf{X} + \lambda \mathbf{I}$$

- “-”: minimise the negative LL (directional curvature is bending upwards)
- check lecture 6 for \mathbf{D}

Example: Polynomial basis regression (28 basis) with l_2 penalty



Example: Radial basis regression (28 basis) with l_2 penalty



l_2 penalty in practice

We do not penalise β_0 (intercept), $I_{p \times p}$ is usually replaced with $I'_{p \times p} = \begin{bmatrix} 0 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ 0 & 0 & \dots & 0 \\ 0 & 0 & \dots & 1 \end{bmatrix}$

!!! Ridge regression is sensitive to the way you scale your data !

- should consider transformation that make the l_2 norm to 1: standardisation (columns of Φ 's norm is 1, i.e. $\|\phi_k\|_2 = 1$)
- other scaling may not produce optimal result
- sk-learn use l_2 penalty by default for logistic regression
 - for numerical stability reason (remember logistic regression's matrix inversion?)

Let's demystify what ridge regression does

$$\beta_{\text{ridge}} = (\Phi^T \Phi + \lambda I)^{-1} \Phi y$$

- obviously, $\lambda \rightarrow 0$, $\beta_{\text{ridge}} \rightarrow \beta_{\text{ml}}$ (which is $\Phi^T \Phi)^{-1} \Phi y$)
- $\lambda \rightarrow \infty$, assume Φ are formed by orthonormal basis (columns), i.e. $\Phi^T \Phi = I_{p \times p}$

$$\beta_{\text{ridge}} = (I + \lambda I)^{-1} \Phi y = \frac{1}{\lambda + 1} \Phi y$$

$\beta \rightarrow \mathbf{0}$ when $\lambda \rightarrow \infty$

- how about those sane λ ?

$$\beta_{\text{ridge}} = \frac{1}{\lambda + 1} \Phi y = \frac{1}{\lambda + 1} \underbrace{(\Phi^T \Phi)^{-1}}_{I^{-1}=I} \Phi y = \underbrace{\frac{1}{\lambda + 1}}_{\text{btw } [0,1]} \beta_{\text{ml}}$$

it shrinks the ML estimator by some percentage.

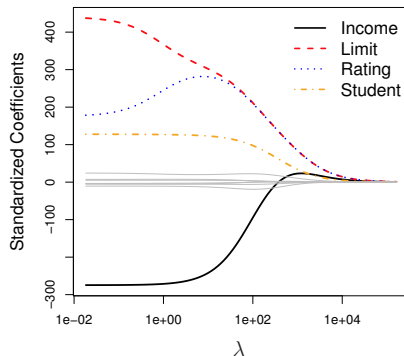
Orthonormal basis are common: e.g. Fourier basis. If Φ is not, the result is still similar: just replace 1 with the eigen value of $\Phi^T \Phi$

Example

$\lambda \rightarrow 0$ (to the left): $\beta_{\text{ridge}} \rightarrow \beta_{\text{ml}}$

$\lambda \rightarrow \infty$ (to the right): $\beta_{\text{ridge}} \rightarrow \mathbf{0}$

other λ : somewhere in between; being shrunk by a percentage



taken from An introduction to statistical learning; Chapter 6, James et al.

Lasso regression l_1 penalty

There is another popular choice for penalty: l_1 norm

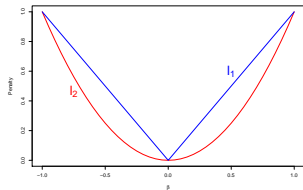
$$L_{\text{lasso}}(\beta) = \underbrace{\sum_{i=1}^m (y^{(i)} - \beta^T \phi(\mathbf{x}^{(i)}))^2}_{\text{previous NLL loss}} + \lambda \underbrace{\|\beta\|_1}_{l_1 \text{ penalty}}$$

- l_1 norm: $\|\mathbf{x}\|_1 = \sum_{j=1}^n |x_j|$, l_2 norm: $\|\mathbf{x}\|_2^2 = \sum_{j=1}^n x_j^2 = \mathbf{x}^T \mathbf{x}$
- ridge shrinks (discount) the ML estimator by some ratio: $\beta_k^{\text{ridge}} = \delta \beta_k^{\text{ml}}$ ($0 < \delta < 1$)
- lasso directly shut them $\beta_k^{\text{lasso}} = 0$
- why though?

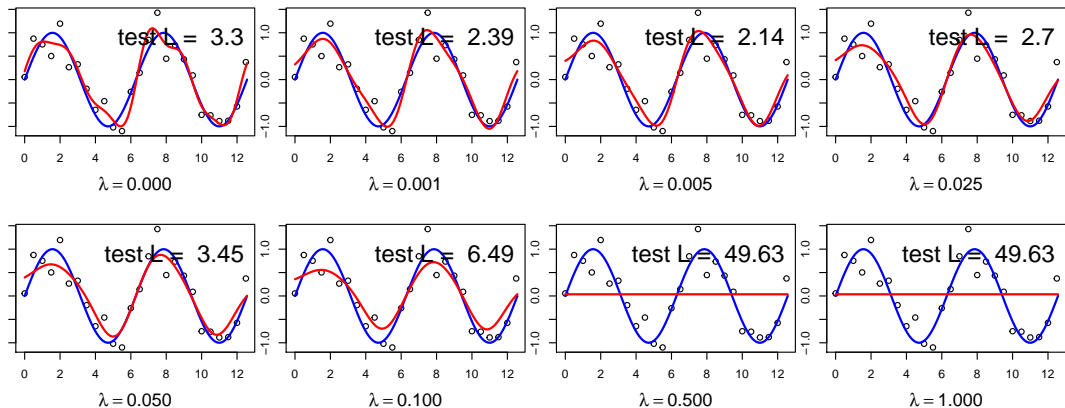
Assume applying gradient descent and $x \rightarrow 0^+$

$$\nabla_x l_1(x) = -1; \quad \nabla_x l_1(x) = -x$$

Penalty l_1 is constantly (strong); but l_2 diminishes



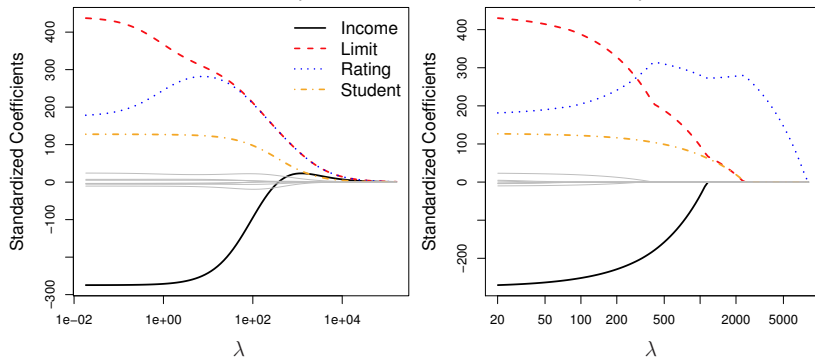
Example: Radial basis regression (28 basis) with l_1 penalty



Compare l_1 and l_2 penalty

Compare and contrast ridge and lasso: (λ has similar effect, i.e. regularisation)

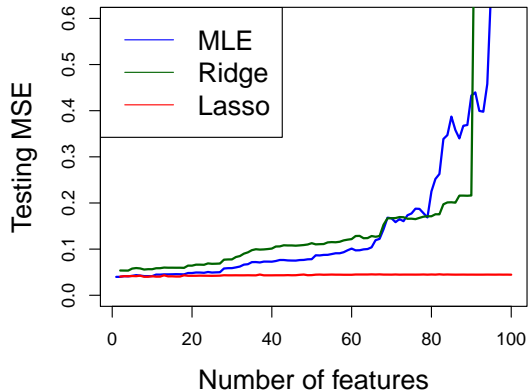
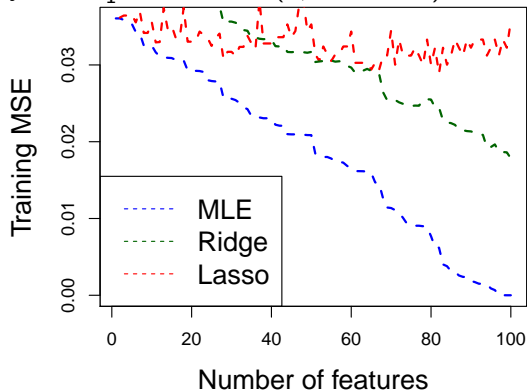
- ridge shrinks the estimator by some percentage (smooth curve)
- Lasso more black and white (feature selection + learning)



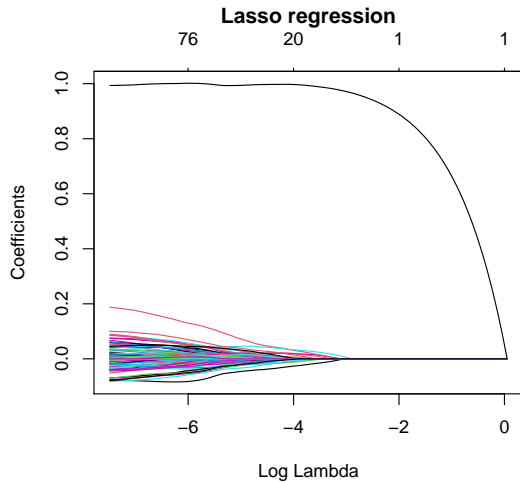
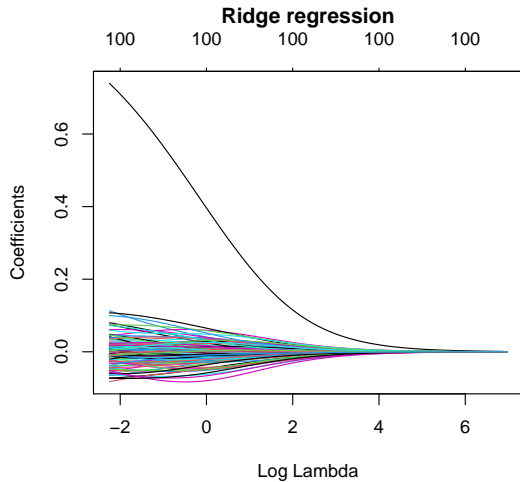
taken from An Introduction to Statistical Learning by James et al.

Compare l_1 and l_2 penalty

Consider \mathbf{X} is a $m = 100$, $n = 100$ (features) matrix sampled from Gaussian ($N(0, 1)$)
 $y^{(i)} = x_1^{(i)} + \epsilon$; $\epsilon \sim N(0, \sigma^2 = 0.2^2)$



The full path



Another view: regularisation is MAP estimation

Regularisation are just Maximum A Posteriori (MAP) estimator

$$P(\beta|\Phi, \mathbf{y}) \propto \underbrace{P(\mathbf{y}|\beta, \Phi)}_{\text{Likelihood}} \underbrace{P(\beta)}_{\text{Prior}}$$

MAP estimation maximise the posterior instead

$$\beta_{\text{MAP}} = \underset{\beta}{\operatorname{argmax}} P(\beta|\Phi, \mathbf{y})$$

Why?

$$\log P(\beta|\Phi, \mathbf{y}) \propto \log P(\mathbf{y}|\beta, \Phi) + \log P(\beta)$$

if we assume $P(\beta) = N(0, \gamma \mathbf{I})$ (a multivariate Gaussian with diagonal covariance)

$$\log P(\beta) = C - \frac{1}{2\gamma} \beta^T \beta$$

Therefore: $\beta_{\text{MAP}} = \beta_{\text{ridge}}$

Suggested reading

- MLAPP 7.5, 13.3, 13.4*
- ESL 3.4
- interesting paper*: A comparison of numerical optimizers for logistic regression by Thomas Minka. <https://tminka.github.io/papers/logreg/minka-logreg.pdf>

Exercise for this lecture (you can discuss it with me in lab session or with your classmates)
ESL 3.12 Show that the ridge regression estimates can be obtained by ordinary least squares regression on an augmented data set. We augment the centered matrix \mathbf{X} with p additional rows $\sqrt{\lambda}\mathbf{I}$, and augment y with p zeros. By introducing artificial data having response zero, the fitting forces a shrinkage.