

Regularization

In data science
flexible = complex \approx many parameters

Q: How do we design **flexible** models which do not overfit the data?

Need to make some assumptions \rightarrow assumptions create bias

No free lunch

But, can still control variance in a highly flexible model as follows: Take polynomial regression

$$Y = \sum_{i=0}^K \beta_j X^j + \epsilon$$

Observation: large K values make model flexible, but also result in very erratic behavior when $K \approx N-1$

Regularization is the idea that we will control this erratic behavior by making β_j small

to implement this we return to optimization picture:

$$\hat{\beta} \text{ minimize } \text{RSS} = \underbrace{\sum_{i=1}^N (y_i - \hat{y}(x_i, D))^2}_{\text{also called loss function}}$$

In regularized regression,

$$\hat{\beta} \text{ minimize } \text{RSS} = \sum_{i=1}^N (y_i - \hat{y}(x_i, D))^2 + R(\hat{\beta})$$

Regularization term which generally increases as $\hat{\beta}$ s increase

In regularized regression we pay a price for larger $\hat{\beta}$

Common Choices

- Ridge regression
(Tikhonov regularization)

$$R(\hat{\beta}) = \lambda \sum_{j=1}^K \beta_j^2$$

controls how much we weight

Note: Later we will see connection between regularization and incorporating prior knowledge in Bayes's statistics

- LASSO (least absolute shrinkage and selection operation)

$$R(\beta) = \lambda \sum_{j=1}^K |\beta_j|$$

Regularization reduces variance b.c. less variation in $\hat{\beta}$ from different data sets
... but increases bias b.c. making $R(\hat{\beta})$ small often means making $y_i - \hat{y}$ larger

Often want to place different weights on different β s, this can be achieved by

$$R(\hat{\beta}) = \sum_{j=1}^K \lambda_j \beta_j^2$$

in python: `OLS-regularized` - see section 2 in Colab

A closer look at regularization and the optimization perspective

Let's consider the case of estimating the mean μ from N samples $\bar{Y}_1, \dots, \bar{Y}_N$

(this is like linear regression with $\beta_0 = \mu, \beta_1 = 0!$)

In this case: $RSS = \sum_{i=1}^N (\bar{Y}_i - \hat{\mu})^2$

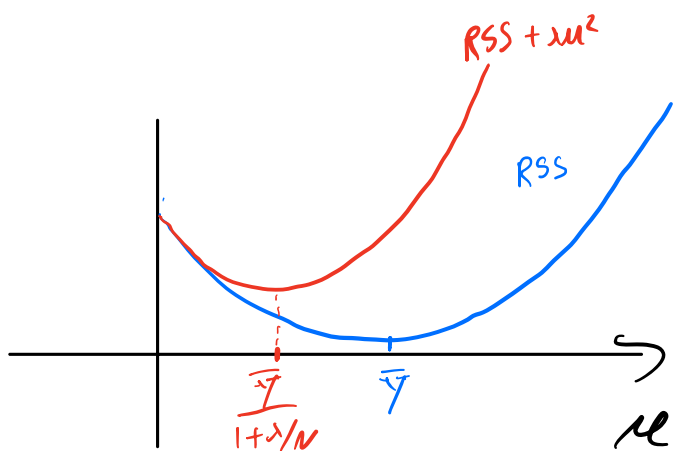
To minimize, solve $\frac{d}{d\mu} RSS = -\sum_{i=1}^N 2(\bar{Y}_i - \hat{\mu}) = 0$

$\Rightarrow \sum_{i=1}^N \bar{Y}_i = N\hat{\mu} \Rightarrow \hat{\mu} = \bar{Y}$ ✓ consistent w/ sample avg./like

Now minimize $L(\mu) = RSS + \lambda\mu^2$

$\Rightarrow \frac{d}{d\mu} L(\mu) = -2\sum_{i=1}^N (\bar{Y}_i - \mu) + 2\lambda\mu = 2N[\bar{Y} - \mu(1 + \lambda/N)] = 0$

$\Rightarrow \hat{\mu} = \frac{\bar{Y}}{1 + \lambda/N}$ larger $\lambda \Rightarrow$ smaller $\hat{\mu}$



more data \Rightarrow regularization matters less

Q: what would happen if we used

$RSS + \lambda(\mu - \mu_0)^2$

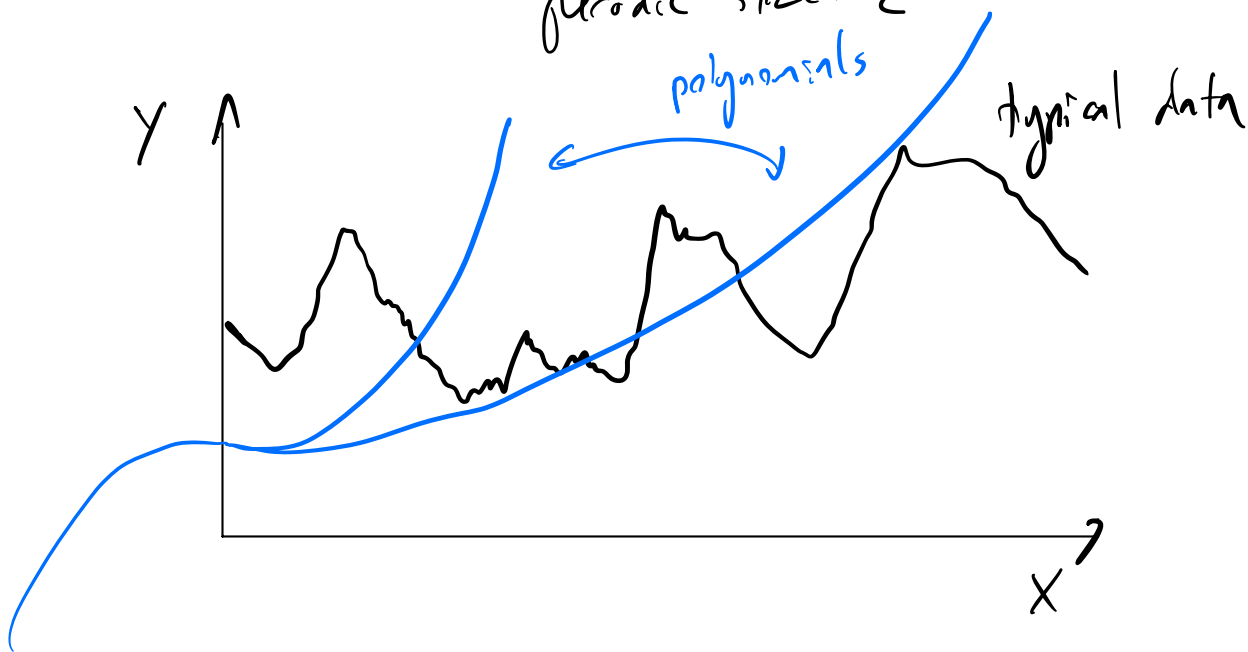
where μ_0 is some fixed \mathbb{R}

Fourier and other models

In general, the monomials x^i are not good terms to use in a regression model, particularly when working with time series data.

This is because: ① x^j and x^i are highly correlated

② $x^j \rightarrow \infty$ as x grows but often we are dealing w/ data which has some approximately periodic structure



In general, when we build a regression model of the form

$$\hat{y} = \sum_{j=0}^K \beta_j \phi_j(x) + \varepsilon$$

we call ϕ_j features or basis functions

Fourier Model

let's consider data where time measurements are uniformly distributed on some interval $[0, L]$:

$$t \sim \text{Uniform}(0, L)$$

here t will play the same role x has been playing

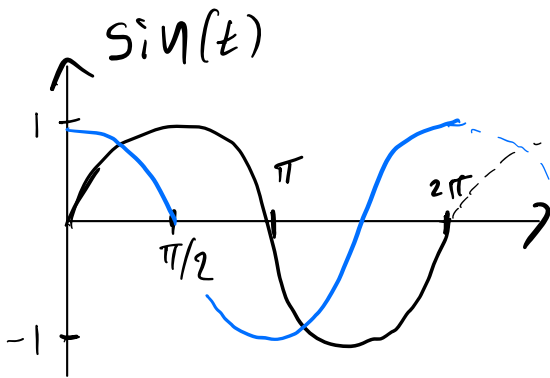
$$Y = \sum_j \beta_j \phi_j(t) + \alpha_j \psi_j(t) + \varepsilon$$

where $\phi_j(t) = \sin(2\pi i x/L)$

$$\psi_j(t) = \cos(2\pi i x/L)$$

note we could rewrite this in usual linear regression form $\sum \beta_j \phi_j$ but it's useful to separate sin and cos terms

Recall:



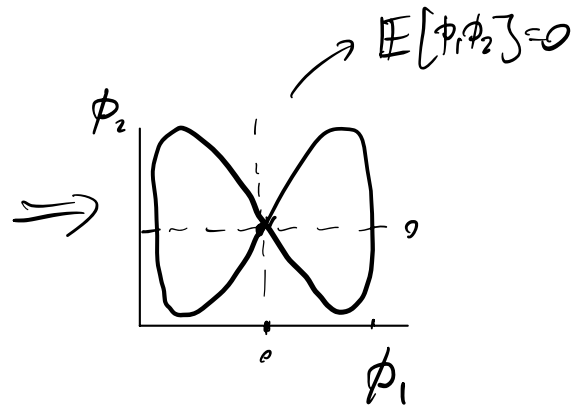
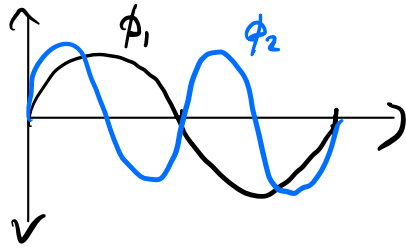
hence ϕ_j, ψ_j are sin/cos functions which go through 1 cycle on $(0, L/i)$

Clearly ϕ_i, ψ_i don't have problem ②
what about problem ①?

Let's look at example:

$$\phi_1 = \sin(2\pi t)$$

$$\phi_2 = \sin(4\pi t)$$



Theorem $t \sim \text{Uniform}(0, 1)$

$$\text{then } \mathbb{E}[\phi_i(t)\phi_j(t)] = 0 \quad i \neq j$$

$$\mathbb{E}[\psi_i(t)\psi_j(t)] = 0 \quad i \neq j$$

$$\mathbb{E}[\gamma_i(t)\gamma_j(t)] = 0 \quad i \neq j$$

note that this means Fourier model is very efficient, each basis function contains different information about our data (different freq.)

this means,

$$\beta_i = \frac{\text{Cov}(Y, \phi_i)}{\text{var}(\phi_i)}, \quad \gamma_i = \frac{\text{Cov}(Y, \psi_i)}{\text{var}(\psi_i)}$$

Warning: in practice there will be some correlation between $\phi_i, \phi_j, \psi_i, \psi_j$ due to randomness in sample

However if t points are not random but evenly spaced these correlations vanish!

See examples in colab notebook