

# Regression and regularization

## Lecture 1d

## Overview

- Linear regression
- Ridge Regression
- Lasso
- Variable selection

732A99/TDDE01

2

1

2

## Simple linear regression

### Model:

$$y \sim N(w_0 + w_1 x, \sigma^2)$$

or

$$y = w_0 + w_1 x + \epsilon, \quad \epsilon \sim N(0, \sigma^2)$$

or

$$p(y|x, w) = N(w_0 + w_1 x, \sigma^2)$$

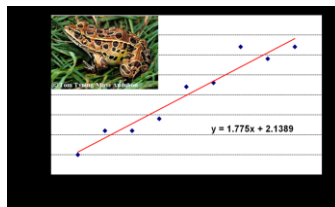
### Terminology:

$w_0$ : intercept (or bias)

$w_1$ : regression coefficient

### Response

The target responds directly and linearly to changes in the feature



732A99/TDDE01

3

3

## Ordinary least squares regression (OLS)

### Model:

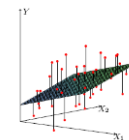
$$y \sim N(w^T x, \sigma^2)$$

where

$$w = \{w_0, \dots, w_d\}$$

$$x = \{1, x_1, \dots, x_d\}$$

Why is "1" here?



The response variable responds directly and linearly to changes in each of the inputs

732A99/TDDE01

4

4

## Ordinary least squares regression

Given data set  $D$

Case	$X_1$	$X_2$		$X_p$	$Y$
1	$x_{11}$	$x_{21}$		$x_{p1}$	$y_1$
2	$x_{12}$	$x_{22}$		$x_{p2}$	$y_2$
3	$x_{13}$	$x_{23}$		$x_{p3}$	$y_3$
$N$	$x_{1N}$	$x_{2N}$		$x_{pN}$	$y_N$

**Estimation:** maximizing the likelihood

$$\hat{w} = \max_w p(D|w)$$

Is equivalent to minimizing

$$RSS(w) = \sum_{i=1}^n (y_i - w^T X_i)^2$$

732A99/TDDE01

5

5

## Matrix formulation of OLS regression

Optimality condition:

$$X^T(y - Xw) = 0$$

where

$$X = \begin{pmatrix} 1 & x_{11} & x_{21} & \dots & x_{p1} \\ 1 & x_{12} & x_{22} & \dots & x_{p2} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{1N} & x_{2N} & \dots & x_{pN} \end{pmatrix} \quad \text{and} \quad y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix}$$

732A99/TDDE01

6

6

## Parameter estimates and predictions

- Least squares estimates of the parameters

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

- Predicted values  $\hat{y} = X\hat{w} = X(X^T X)^{-1} X^T y = Py$

- Linear regression belongs to the class of **linear smoothers**



Hat matrix

Why is it called so?

732A99/TDDE01

7

7

## Degrees of freedom

Definition:

$$df(\hat{y}) = \frac{1}{\sigma^2} \sum_{i=1}^N \text{Cov}(\hat{y}_i, y_i)$$

- Larger covariance  $\rightarrow$  stronger connection  $\rightarrow$  model can approximate data better  $\rightarrow$  model more flexible (complex)
- For linear smoothers  $\hat{Y} = S(X)Y$

$$df = \text{trace}(S)$$

- For linear regression, degrees of freedom is  $df = \text{trace}(P) = p$

732A99/TDDE01

8

8

## Different types of features

- Interval variables
- Numerically coded ordinal variables
  - (small=1, medium=2, large=3)
- Dummy coded qualitative variables

**Example of dummy coding:**

$$x_1 = \begin{cases} 1, & \text{if Jan} \\ 0, & \text{otherwise} \end{cases}$$

$$x_2 = \begin{cases} 1, & \text{if Feb} \\ 0, & \text{otherwise} \end{cases}$$

.

.

$$x_{11} = \begin{cases} 1, & \text{if Nov} \\ 0, & \text{otherwise} \end{cases}$$

**Basis function expansion:**

If  $y = w_0 + w_1 x_1 + w_2 x_1^2 + w_3 e^{-x_2} + \epsilon$ ,

Model becomes linear if to recompute:

$$\begin{aligned} \phi_1(x_1) &= x_1 \\ \phi_2(x_1) &= x_1^2 \\ \phi_3(x_1) &= e^{-x_2} \end{aligned}$$

732A99/TDDE01

9

9

## Basis function expansion

- In general  $\phi_1(\dots)$  may be a function of several  $x$  components
- Having data given by  $\mathbf{X}$ , compute new data

$$\Phi = \begin{pmatrix} 1 & \phi_1(x_{11}, \dots, x_{1p}) & \dots & \phi_p(x_{11}, \dots, x_{1p}) \\ \vdots & \vdots & \ddots & \vdots \\ 1 & \phi_1(x_{n1}, \dots, x_{np}) & \dots & \phi_p(x_{n1}, \dots, x_{np}) \end{pmatrix}$$

- If doing a basis function in a model, replace  $\mathbf{X}$  by  $\Phi$  everywhere where  $\mathbf{X}$  is used:

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

732A99/TDDE01

10

10

## Linear regression in R

- `fit=lm(formula, data, subset, weights,...)`
  - data** is the data frame containing the predictors and response values
  - formula** is expression for the model
  - subset** which observations to use (training data)?
  - weights** should weights be used?

**fit** is object of class **lm** containing various regression results.

- Useful functions (many are generic, used in many other models)
  - Get details about the particular function by `"?",` for ex. `predict.lm`

```
summary(fit)
predict(fit, newdata, se.fit, interval)
coefficients(fit) # model coefficients
confint(fit, level=0.95) # CIs for model parameters
fitted(fit) # predicted values
residuals(fit) # residuals
```

732A99/TDDE01

11

11

## An example of ordinary least squares regression

```
mydata=read.csv2("Bilexempel.csv")
fit1=lm(Price~Year, data=mydata)
summary(fit1)
fit2=lm(Price~Year+Mileage+Equipment,
data=mydata)
summary(fit2)
```

**Response variable:**  
Requested price of used Porsche cars  
(1000 SEK)

```
> summary(fit1)
Call:
lm(formula = Price ~ Year, data = mydata)

Residuals:
    min       1Q   median       3Q      max
-167683  -14683   20056  35933  72317

Coefficients:
(Intercept) 78161027 8448038 -9.232 6.00e-13 ***
Year        38246     4226  9.288 5.25e-13 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 37270 on 57 degrees of freedom
Multiple R-squared:  0.5932
F-statistic: 86.26 on 1 and 57 DF,  p-value: 5.248e-13
```

**Inputs:**  
 $X_1$  = Manufacturing year  
 $X_2$  = Mileage (km)  
 $X_4$  = Equipment (0 or 1)

732A99/TDDE01

12

12

## An example of ordinary least squares regression

```
> summary(fit2)

Call:
lm(formula = Price ~ Year + Mileage + Equipment, data = mydata)

Residuals:
    Min       1Q   Median       3Q      Max
-66223 -10525    -739   14128   65332

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) -2.081e+07  6.309e+06  -3.302  0.00169 **
Year          1.062e+04  3.154e+03   3.366  0.00139 **
Mileage      -2.077e+00  5.022e-01 -0.269  0.78414 ***
Equipment     5.790e+04  1.041e+04   5.563  8.08e-07 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 29270 on 55 degrees of freedom
Multiple R-squared:  0.8997, Adjusted R-squared:  0.8942
F-statistic: 164.5 on 3 and 55 DF,  p-value: < 2.2e-16
```

732A99/TDDE01

13

13

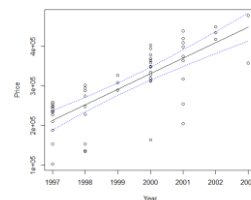
## An example of ordinary least squares regression

### • Prediction

```
fitted <- predict(fit1, interval =
"confidence")

# plot the data and the fitted line
attach(mydata)
plot(Year, Price)
lines(Year, fitted[, "fit"])

# plot the confidence bands
lines(Year, fitted[, "lwr"], lty = "dotted",
col="blue")
lines(Year, fitted[, "upr"], lty = "dotted",
col="blue")
detach(mydata)
```



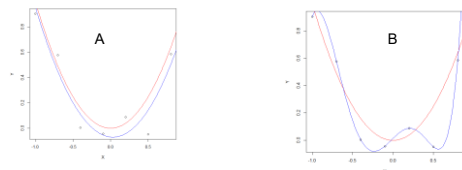
732A99/TDDE01

14

14

## Ridge regression

- Problem: linear regression can overfit:
  - Take  $Y := Y, X_1 = X, X_2 = X^2, \dots, X_p = X^p \rightarrow$  polynomial model, fit by linear regression
  - High degree of polynomial leads to overfitting.



732A99/TDDE01

15

15

## Ridge regression

- Idea: Keep all predictors but shrink coefficients to make model less complex  
 minimize  $-\log\text{likelihood} + \lambda_0 \|w\|_2^2$   
 $\rightarrow$   **$l_2$  regularization**
  - Given that model is Gaussian, we get **Ridge regression**:

$$\hat{w}^{\text{ridge}} = \underset{w}{\operatorname{argmin}} \left\{ \sum_{i=1}^N (y_i - w_0 - w_1 x_{i1} - \dots - w_p x_{ip})^2 + \lambda \sum_{j=1}^p w_j^2 \right\}$$

- $\lambda > 0$  is **penalty factor**

732A99/TDDE01

16

16

## Ridge regression

Equivalent form

$$\hat{w}^{ridge} = \operatorname{argmin}_{\mathbf{w}} \sum_{i=1}^N (y_i - w_0 - w_1 x_{1i} - \dots - w_p x_{pi})^2$$

subject to  $\sum_{j=1}^p w_j^2 \leq s$

**Solution**

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

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

Hat matrix

How do we  
compute degrees  
of freedom here?

732A99/TDDE01

17

17

## Ridge regression

### Properties

- Extreme cases:
  - $\lambda = 0$  usual linear regression (no shrinkage)
  - $\lambda = +\infty$  fitting a constant ( $w = 0$  except of  $w_0$ )
- When input variables are orthogonal (not realistic),  $\mathbf{X}^T \mathbf{X} = \mathbf{I} \rightarrow$ 

$$\hat{\mathbf{w}}^{ridge} = \frac{1}{1+\lambda} \mathbf{w}^{linreg} \rightarrow \text{coefficients are equally shrunk}$$
- **Ridge regression is particularly useful if the explanatory variables are strongly correlated to each other.**
  - Correlated variables often correspond large  $w \rightarrow$  shrunk
- Degrees of freedom decrease when  $\lambda$  increases
  - $\lambda = 0 \rightarrow d.f. = p$

732A99/TDDE01

18

18

## Ridge regression

### Properties

- Shrinking enables estimation of regression coefficients even if the number of parameters exceeds the number of cases! ( $\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}$  is always nonsingular)
  - Compare with linear regression
- How to estimate  $\lambda$ ?
  - cross-validation

732A99/TDDE01

19

19

## Ridge regression

- **Bayesian view**
  - Ridge regression is just a special form of Bayesian Linear Regression with constant  $\sigma^2$ :

$$\mathbf{y} \sim \mathcal{N}(\mathbf{y} | \mathbf{w}_0 + \mathbf{X} \mathbf{w}, \sigma^2 \mathbf{I})$$

$$\mathbf{w} \sim \mathcal{N}\left(\mathbf{0}, \frac{\sigma^2}{\lambda} \mathbf{I}\right)$$

**Theorem** MAP estimate to the Bayesian Ridge is equal to solution in frequentist Ridge

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

- In Bayesian version, we can also make inference about  $\lambda$

732A99/TDDE01

20

20

## Ridge regression

**Example Computer Hardware Data Set** : performance measured for various processors and also

- Cycle time
- Memory
- Channels
- ...

Build model predicting performance



732A99/TDDE01

21

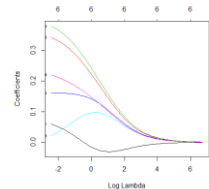
21

## Ridge regression

- R code: use package **glmnet** with  $\alpha=0$  (Ridge regression)
- Seeing how Ridge converges

```
data=read.csv("machine.csv", header=F)
covariates=scale(data[,3:8])
response=scale(data[, 9])

model0=glmnet(as.matrix(covariates),
response, alpha=0,family="gaussian")
plot(model0, xvar="lambda", label=TRUE)
```



732A99/TDDE01

22

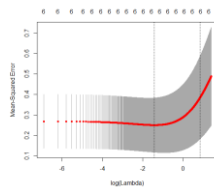
22

## Ridge regression

- Choosing the best model by cross-validation:

```
model=cv.glmnet(as.matrix(covariates),
response, alpha=0,family="gaussian")
model$lambda.min
plot(model)
coef(model, s="lambda.min")
```

```
> coef(model, s="lambda.min")
7 x 1 sparse Matrix of class "dgCMatrix"
(Intercept) -4.530442e-17
v3          3.420739e-02
v4          3.085696e-01
v5          3.403839e-01
v6          1.593470e-01
v7          5.489116e-02
v8          1.970982e-01
```



```
> model$lambda.min
[1] 0.046
```

732A99/TDDE01

23

23

## Ridge regression

- How good is this model in prediction?

```
ind=sample(209, floor(209*0.5))
data1=scale(data[,3:9])
train=data1[ind,]
test=data1[-ind,]
```

```
covariates=train[,1:6]
response=train[, 7]
model=cv.glmnet(as.matrix(covariates), response, alpha=1, family="gaussian",
lambda=seq(0,1,0.001))
y=test[,7]
ynew=predict(model, newx=as.matrix(test[, 1:6]), type="response")

#Coefficient of determination
sum((ynew-mean(y))^2)/sum((y-mean(y))^2)
```

Note that data are so small so numbers change much for other train/test

```
sum((ynew-y)^2)
```

```
> sum((ynew-mean(y))^2)/sum((y-mean(y))^2)
[1] 0.5438148
> sum((ynew-y)^2)
[1] 18.04988
> 1
```

732A99/TDDE01

24

24

## LASSO

- **Idea:** Similar idea to Ridge
- Minimize minus loglikelihood plus **linear** penalty factor  $\rightarrow I_1$  regularization
  - Given that model is Gaussian, we get **LASSO** (least absolute shrinkage and selection operator):

$$\hat{w}^{lasso} = \operatorname{argmin} \left\{ \sum_{i=1}^N (y_i - w_0 - w_1 x_{1i} - \dots - w_p x_{pi})^2 + \lambda \sum_{j=1}^p |w_j| \right\}$$

- $\lambda > 0$  is **penalty factor**



732A99/TDDE01

25

25

## LASSO

- Equivalently

$$\hat{w}^{lasso} = \operatorname{argmin} \sum_{i=1}^N (y_i - w_0 - w_1 x_{1i} - \dots - w_p x_{pi})^2$$

subject to  $\sum_{j=1}^p |w_j| \leq s$

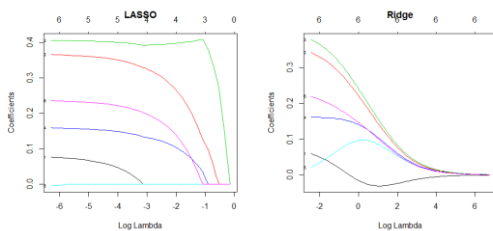
732A99/TDDE01

26

26

## LASSO vs Ridge

- **LASSO yields sparse solutions!**
- Example** Computer hardware data



732A99/TDDE01

27

27

## LASSO vs Ridge

- Only 5 variables selected by LASSO

```
> coef(model1, s="lambda.min")
7 x 1 sparse Matrix of class "dgCMatrix"
      1
(Intercept) -5.091825e-17
v3          6.350488e-02
v4          3.578607e-01
v5          4.033670e-01
v6          1.541329e-01
v7          .
v8          2.287134e-01
~ 1
```

$> \text{sum}((y_{\text{new}} - \text{mean}(y))^2) / \text{sum}((y - \text{mean}(y))^2)$   
 $[1] 0.5826904$   
 $> \text{sum}((y_{\text{new}} - y)^2)$   
 $[1] 16.63756$

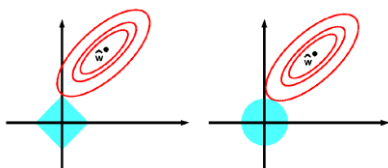
732A99/TDDE01

28

28

## LASSO vs Ridge

- Why Lasso leads to sparse solutions?
  - Feasible area for Ridge is a circle (2D)
  - Feasible area for LASSO is a polygon (2D)



732A99/TDDE01

29

29

## LASSO properties

- Lasso is widely used when  $p \gg n$** 
  - Linear regression breaks down when  $p > n$
  - Application: DNA sequence analysis, Text Prediction

- When inputs are orthonormal,

$$\hat{w}_i^{\text{lasso}} = \text{sign}(w_i^{\text{linreg}}) \left( |w_i^{\text{linreg}}| - \frac{\lambda}{2} \right)_+$$

- No explicit formula for  $\hat{w}^{\text{lasso}}$ 
  - Optimization algorithms used

**Coding in R: use**  
**glmnet() with**  
**alpha=1**

732A99/TDDE01

30

30

## Variable selection

- .. Or "Feature selection"

Often, we do not need all features available in the data to be in the model

### Reasons:

- Model can become overfitted (recall polynomial regression)
- Large number of predictors → model is difficult to use and interpret

732A99/TDDE01

31

31

## Variable selection

### Alternative 1: Variable subset selection

- Best subset selection:**
  - Consider different subsets of the full set of features, fit models and evaluate their quality
    - Problem: computationally difficult for  $p$  around 30 or more
    - How to choose the best model size? Some measure of predictive performance normally used (ex. AIC).
- Forward and Backward stepwise selection**
  - Starts with 0 features (or full set) and then adds a feature (removes feature) that most improves the measure selected.
    - Can handle large  $p$  quickly
    - Does not examine all possible subsets (not the "best")

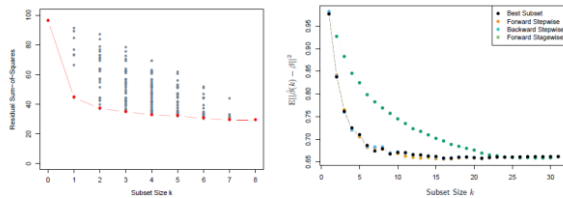
732A99/TDDE01

32

32



## RSS and MSE depend on k



732A99/TDDE01

33

## Variable selection in R

- Use `stepAIC()` in MASS

```
library(MASS)
fit <- lm(V9~., data=data.frame(data1))
step <- stepAIC(fit, direction="both")
step$anova
summary(step)
```

```
[Call:
lmFormula = V9 ~ V3 + V4 + V5 + V6 + V8, data = data.fry]

Residuals:
    Min       1Q   Median       3Q      Max
-1.30232  -0.33532   0.03378   0.30507   2.42280

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) -5.780e-17  2.536e-02   0.000  1.0000
V3           7.944e-02  2.820e-02   2.813  0.0054 **
V4           3.462e-01  4.320e-02   8.000 4.36e-15 ***
V5           4.051e-01  4.684e-02   8.635 1.18e-15 ***
V6           1.354e-01  3.394e-02   4.007 5.07e-06 ***
V8           2.300e-01  3.350e-02   7.031 3.06e-11 ***
```

```
> step <- stepAIC(fit, direction="both")
Start: AIC=-405.35
V9 ~ V3 + V4 + V5 + V6 + V7 + V8

              Df Sum of Sq  RSS   AIC
<none>                 28 117  -407.25
+ V7                   1  0.0139 28.117  -407.25
+ V3                   1  1.0819 29.185  -399.46
+ V6                   1  2.9383 31.041  -386.57
+ V8                   1  6.3150 34.418  -364.99
+ V4                   1  9.7492 37.852  -345.11
+ V5                   1 10.4837 38.586  -341.09

Step: AIC=-407.25
V9 ~ V3 + V4 + V5 + V6 + V8

              Df Sum of Sq  RSS   AIC
<none>                 28 117  -407.25
+ V7                   1  0.0139 28.103  -405.35
+ V3                   1  1.0958 29.212  -401.26
+ V6                   1  3.0431 31.160  -387.77
+ V8                   1  6.8472 34.964  -363.70
+ V4                   1  9.8840 38.101  -345.74
+ V5                   1 10.4713 38.586  -343.08
```

732A99/TDDE01

34