# Administration

- ▶ Homework 1 available Thursday
- ▶ Discussion of project requirements on Thursday
- ▶ NSERC summer undergraduate awards
- ▶ Fields-MITACS undergraduate summer research

  http://www.fields.utoronto.ca/programs/scientific/10-11/summer-research/

# Geometric view of least squares fitting

- $\hat{\beta} = (X^T X)^{-1} X^T y$
- $\hat{\beta} = (\hat{\beta}_0, \hat{\beta}_1, \ldots, \hat{\beta}_p)$
- $\hat{\beta}_p$ can be obtained by a series of regressions (projections) as outlined in algorithm 3.1 on p.54

   regress $\mathbf{x}_1$ on 1, get coefficient $\hat{\gamma}_{01}$, form residual $z_1 = \mathbf{x}_1 - \hat{\mathbf{x}}_1$
   regress $\mathbf{x}_2$ on 1, $z_1$, get coefs $\hat{\gamma}_{02}, \hat{\gamma}_{12}$, form residual $z_2 = \mathbf{x}_2 - \hat{\gamma}_{02}1 - \hat{\gamma}_{12}z_1$
   $\vdots$
   regress $\mathbf{x}_p$ on $z_{p-1}, z_{p-2}, \ldots, z_1, 1$ to get $z_p = \mathbf{x}_p - \hat{\mathbf{x}}_p$
   regress $y$ on $z_p$ to get $\hat{\beta}_p$

- illustration on prostate training data – see
  `prostateRsession.txt`

# QR Decomposition $X = Z\Gamma$

- matrix representation $X = Z\Gamma$
- $Z$ has columns $z_j$
-
$$\Gamma = \begin{pmatrix} 0 & \hat{\gamma}_{01} & \hat{\gamma}_{02} & \dots & \hat{\gamma}_{0p} \\ 0 & 0 & \hat{\gamma}_{12} & \dots & \hat{\gamma}_{1p} \\ & & \ddots & & \\ & & & & \hat{\gamma}_{p-1,p} \\ 0 & \dots & & & 0 \end{pmatrix}$$

- Let $D_{jj} = (z_j^T z_j)^{1/2} = ||z_j||$ and $D = \text{diag}(D_{jj})$ dimension?
- $X = ZD^{-1}D\Gamma = QR$
- $Q^T Q = I, \quad R$ is upper triangular
- $\hat{\beta} = R^{-1}Q^T y$
- $\hat{y} = QQ^T y$ check `pr.lm$qr`
- Gram-Schmidt

3 / 24

# **Singular value decomposition,** $X = UDV$

- $X$ assumed to be centered, so all columns add to zero (p.64, l.-4)
- $U$ is $N \times p$, $V$ is $p \times p$
- $D$ is diagonal, $d_1 \geq d_2 \geq \cdots \geq d_p \geq 0$    not the same $D$
- 

$$X^T X = V D^2 V^T$$

  eigendecomposition of $X^T X$ and of $NS = XX^T$
- $\hat{y} = X\hat{\beta} = UU^T y$
    $= QQ^T y$ different orthogonal bases for ...
- define

$$z_1 = Xv_1 = u_1 d_1$$

- note that $\mathrm{var}(z_1) = d_1^2 / N$
- $z_1$ is the derived variable with the largest variance: the first principal component of $X$
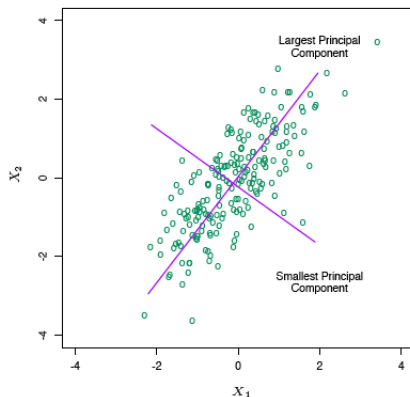- $z_2$ has the largest variance among linear combinations orthogonal to $z_1$

**FIGURE 3.9.** *Principal components of some input data points. The largest principal component is the direction that maximizes the variance of the projected data, and the smallest principal component minimizes that variance. Ridge regression projects* **y** *onto these components, and then shrinks the coefficients of the low–variance components more than the high-variance components.*

# ... singular value decomposition

►

$$X_{N\times p} = U_{N\times p}D_{p\times p}V_{p\times p}^T, \quad U^T U = I, \quad V^T V = I,$$
$$D = \operatorname{diag}(d_1, \ldots d_p)$$

$$
\begin{aligned}
X\hat{\beta}_{LS} &= X(X^T X)^{-1}X^T y \\
&= UDV^T(VD{\color{red}U^T U}DV^T)^{-1}VDU^T y \\
&= UD{\color{red}V^T V}^{T^{-1}}D^{-2}{\color{red}V^{-1} V}DU^T y \\
&= UU^T y = \Sigma_{j=1}^p u_j u_j^T y
\end{aligned}
$$

`svd(model.matrix(pr.lm))`, for example

# Model selection: subsets (§3.3.1)

- linear regression: forward, backward, stepwise, all possible subsets regression

- $RSS(\hat{\beta}) = (y - X\hat{\beta})^T(y - X\hat{\beta})$ $\qquad = \Sigma(y_i - \hat{y}_i)^2$
  Figure 3.5 $\qquad\qquad$ <sub>this is called SSE in the 302 text</sub>

- if we add a regressor, say from $X_{p-1}$ to $X_p$, $RSS(\hat{\beta})$ necessarily decreases

- forward (stepwise) selection starts with one predictor (usually the constant term) and stops when no additional predictor is statistically significant `step(pr.lm, direction = "forward", ...)`

- backward (stepwise) selection starts with all predictors and deletes least significant ... `direction = "backward"`...

- stepwise selection checks at each stage whether or not to add variables back in `direction = "both"`

# ... subset selection

- ► forward stagewise: a "slow" version of forward stepwise, in which coefficients are not re-computed
- ► all possible subsets regression considers all $2^p$ models.
- ► for $p < 30$, feasible with the "leaps and bounds" algorithm, implemented in package `leaps` (See Figure 3.6), also `regsubsets`
- ► Figure 3.7
- ► huh?
- ► 10-fold cross-validation
- ► model selection related to expected prediction error: theory to come in Ch. 7

# "Mallows' $C_p$"

▶ a common adjustment to measure benefit of adding further parameters:

$$C_p = \frac{RSS_p}{\sigma^2} + 2p - N \quad \text{if } \sigma^2 \text{ is known;}$$

▶ or an estimate of this, if $\sigma^2$ is unknown
▶ rule of thumb: choose $p$ so that $C_p$ is small and $C_p \simeq p$
▶ can be shown to be a good choice for prediction (details deferred until Chapter 7)
▶ a closely related, more general, criterion *AIC* (Akaike's Information Criterion)
▶ for our linear model

$$AIC \approx N \log(RSS_p/N) + 2p + \textit{constant}$$

▶ stepAIC in MASS library; step in base

**prostate data: all possible subsets**
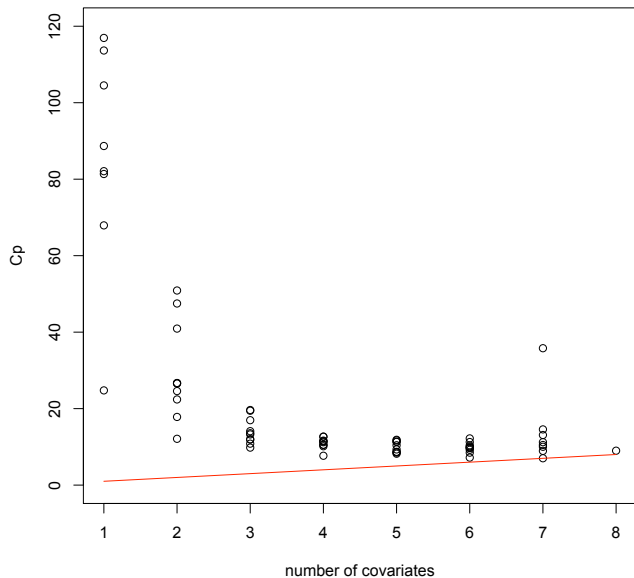
Cp vs. number of covariates

Figure 3.5

# step

```
> step(pr.lm)
Start:  AIC=-37.13
lpsa ~ lcavol + lweight + age + lbph + svi + lcp + gleason +
    pgg45

          Df Sum of Sq    RSS     AIC
- gleason  1     0.011  29.437 -39.103
<none>                  29.426 -37.128
- age      1     0.989  30.415 -36.914
- pgg45    1     1.532  30.959 -35.727
- lcp      1     1.768  31.195 -35.218
- lbph     1     2.144  31.571 -34.415
- svi      1     3.093  32.520 -32.430
- lweight  1     3.839  33.265 -30.912
- lcavol   1    14.610  44.037 -12.118

Step:  AIC=-39.1
lpsa ~ lcavol + lweight + age + lbph + svi + lcp + pgg45

          Df Sum of Sq    RSS     AIC
<none>                  29.437 -39.103
- age      1     1.102  30.540 -38.639
- lcp      1     1.758  31.196 -37.216
- lbph     1     2.135  31.573 -36.411
- pgg45    1     2.376  31.813 -35.903
- svi      1     3.166  32.604 -34.258
- lweight  1     4.005  33.442 -32.557
- lcavol   1    14.887  44.325 -13.681

Call:
lm(formula = lpsa ~ lcavol + lweight + age + lbph + svi + lcp +      pgg45, data = train)
```

# Shrinkage Methods: (§3.4)

▶ Ridge regression

▶

$$\hat{\beta}_{LS} = (X^T X)^{-1} X^T y$$
$$\hat{\beta}_{ridge} = (X^T X + \lambda I)^{-1} X^T y$$

▶ can show that $\hat{\beta}_{ridge}$ satisfies

$$\min_{\beta} \left( \Sigma \{ y_i - \beta_0 - \Sigma_{j=1}^{p} x_{ij} \beta_j \}^2 + \lambda \Sigma_{j=1}^{p} \beta_j^2 \right)$$

$$\min_{\beta} \Sigma \{ y_i - \beta_0 - \Sigma_{j=1}^{p} x_{ij} \beta_j \}^2 \quad \text{s.t. } \Sigma \beta_j^2 \leq t$$

▶ Assume $x_j$'s are centered and put these in matrix $X$ (with no column of 1's:

$$\min_{\beta} (y - X\beta)^T (y - X\beta) \quad \text{s.t. } ||\beta||^2 \leq t$$

# ... ridge regression

- $\min_\beta \{(y - X\beta)^T(y - X\beta) + \lambda||\beta||^2\}$
- $\lambda$ is a tuning parameter: $\lambda = 0$ gives $\hat{\beta}_{LS}$, $\lambda \to \infty$

  Figure 3.8

- in R the library MASS library(MASS ) has a ridge regression version of lm called lm.ridge
- if columns of $X$ are nearly linearly dependent (multicollinearity), $\hat{\beta}$'s for these columns should be shrunk towards 0.
- essential that the predictors are all scaled to the same units
- this is difficult for interpretation of the coefficients

$$
\begin{aligned}
X\hat{\beta}_{ridge} &= X(X^T X + \lambda I)^{-1} X^T y \\
&= UDV^T (VD^2 V^T + \lambda I)^{-1} VDU^T y \\
&= UDV^T (VD^2 V^T + \lambda VV^T)^{-1} VDU^T y \\
&= UD(D^2 + \lambda I)^{-1} DU^T y \\
&= \Sigma_{j=1}^{p} u_j \left(\frac{d_j^2}{d_j^2 + \lambda}\right) u_j^T y
\end{aligned}
$$

$$
df(\lambda) = \mathrm{tr}[X(X^T X + \lambda I)^{-1} X^T] = \Sigma_{j=1}^{p} \frac{d_j^2}{d_j^2 + \lambda}
$$

Figure 3.7

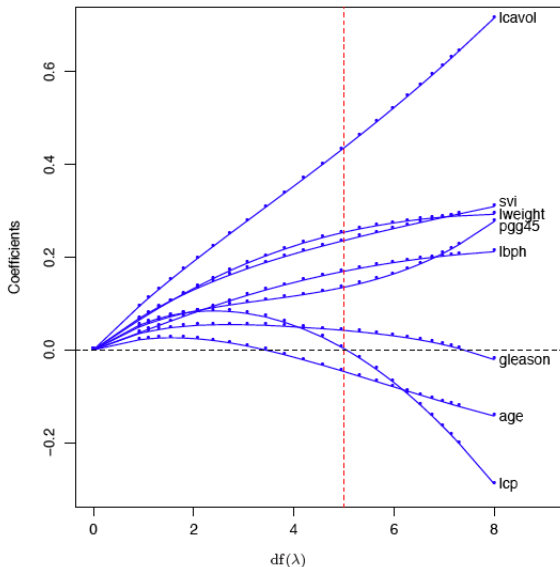$df(\lambda)$ called effective number of parameters in Ch. 7

**FIGURE 3.8.** *Profiles of ridge coefficients for the prostate cancer example, as the tuning parameter $\lambda$ is varied. Coefficients are plotted versus $\mathrm{df}(\lambda)$, the effective degrees of freedom. A vertical line is drawn at $\mathrm{df} = 5.0$, the value chosen by cross-validation.*

# Lasso

- $$\min_\beta \left( \Sigma\{y_i - \beta_0 - \Sigma_{j=1}^p x_{ij}\beta_j\}^2 + \lambda\Sigma_{j=1}^p |\beta_j| \right)$$

- $$\min_\beta \Sigma\{y_i - \beta_0 - \Sigma_{j=1}^p x_{ij}\beta_j\}^2 \quad \text{s.t. } \Sigma|\beta_j| \le t$$

- ▶ quadratic programming problem
- ▶ $\hat{\beta}^{lasso}$ is nonlinear function of $y$
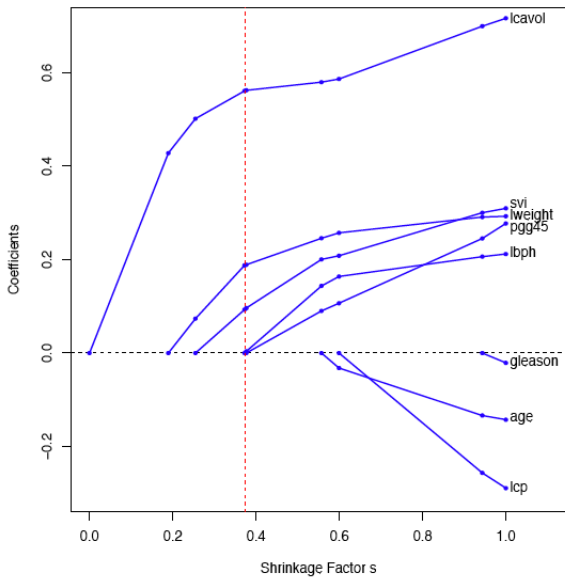- ▶ Figure 3.10
- ▶ Table 3.3

**FIGURE 3.10.** *Profiles of lasso coefficients, as the tuning parameter $t$ is varied. Coefficients are plotted versus $s = t / \sum_1^p |\hat{\beta}_j|$. A vertical line is drawn at $s = 0.36$, the value chosen by cross-validation. Compare Figure 3.8 on page 65; the lasso*

**TABLE 3.3.** *Estimated coefficients and test error results, for different subset and shrinkage methods applied to the prostate data. The blank entries correspond to variables omitted.*

| Term | LS | Best Subset | Ridge | Lasso | PCR | PLS |
|---|---|---|---|---|---|---|
| Intercept | 2.465 | 2.477 | 2.452 | 2.468 | 2.497 | 2.452 |
| lcavol | 0.680 | 0.740 | 0.420 | 0.533 | 0.543 | 0.419 |
| lweight | 0.263 | 0.316 | 0.238 | 0.169 | 0.289 | 0.344 |
| age | −0.141 | | −0.046 | | −0.152 | −0.026 |
| lbph | 0.210 | | 0.162 | 0.002 | 0.214 | 0.220 |
| svi | 0.305 | | 0.227 | 0.094 | 0.315 | 0.243 |
| lcp | −0.288 | | 0.000 | | −0.051 | 0.079 |
| gleason | −0.021 | | 0.040 | | 0.232 | 0.011 |
| pgg45 | 0.267 | | 0.133 | | −0.056 | 0.084 |
| Test Error | 0.521 | 0.492 | 0.492 | 0.479 | 0.449 | 0.528 |
| Std Error | 0.179 | 0.143 | 0.165 | 0.164 | 0.105 | 0.152 |

# ...Lasso

- ▶ in Table 3.3 each method had a tuning parameter to choose; they used cross-validation within the training data
- ▶ in `lm.ridge` you can extract a component called `$GCV`
- ▶ the quantity $\Sigma d_j^2/(d_j^2 + \lambda)$ has an interpretation as the number of 'degrees of freedom' or number of 'parameters' used by the ridge regression fit
- ▶ book says that the best value is 4.16, which corresponds to quite a large $\lambda$ (39); the GCV criterion chooses $\lambda = 5$
- ▶ analysis of lasso more difficult; note Figure 3.10 plotted against $t/\Sigma|\hat{\beta}_j|$

# ... smoothing

- ▶ ridge regression gives "proportional shrinkage"
- ▶ subset selection gives "hard thresholding" (some $\beta_j \to 0$)
- ▶ lasso gives "soft thresholding": blend of shrinkage and zeroing (Figure 3.10 and Figure 3.11)
- ▶ Least Angle Regression (LAR): combine forward stagewise regression with the lasso
- ▶ related to the Dantzig selector (Candes and Tao, AS 2007)

# Mean squared error of prediction in linear models

Let $\tilde{\beta} = \tilde{\beta}(y)$ be a competing estimator of $\beta$ (not $\hat{\beta}$, the LS estimator). Using $\tilde{\beta}$ for prediction would give $\tilde{y}_0 = x_0^T \tilde{\beta}$, where $x_0^T = (1, x_{01}, \ldots, x_{0p})$ is the new value of the inputs. The expected prediction error is

$$
\begin{aligned}
E(\tilde{y}_0 - y_0)^2 &= E(x_0^T \tilde{\beta} - y_0)^2 \quad \text{(by definition)} \\
&= E(y_0 - x_0^T \beta + x_0^T \beta - x_0^T \tilde{\beta})^2 \\
&= \text{var}(y_0) + E(x_0^T \tilde{\beta} - x_0^T \beta)^2 \quad \text{(why is cross prod 0?)} \\
&= \sigma^2 + E\{x_0^T(\tilde{\beta} - \beta)(\tilde{\beta} - \beta)^T x_0\} \\
&= \sigma^2 + x_0^T E\{(\tilde{\beta} - \beta)(\tilde{\beta} - \beta)^T\} x_0 \\
&= \sigma^2 + x_0^T E\{(\tilde{\beta} - E\tilde{\beta} + E\tilde{\beta} - \beta)(\tilde{\beta} - E\tilde{\beta} + E\tilde{\beta} - \beta)x_0 \\
&= \sigma^2 + x_0^T [E\{(\tilde{\beta} - E\tilde{\beta})(\tilde{\beta} - E\tilde{\beta})^T\} + (E\tilde{\beta} - \beta)(E\tilde{\beta} - \beta)^T] x_0 \\
&= \sigma^2 + x_0^T \{\text{cov}(\tilde{\beta}) + \text{bias}^2(\tilde{\beta})\} x_0
\end{aligned}
$$

The first term, $\sigma^2$, is unavoidable. The next two terms together are the Mean Squared Error (MSE) of the prediction $\tilde{y}_0$, and are shown here to be a function of $x_0$ and the MSE of $\tilde{\beta}$. If $\tilde{\beta}$ is unbiased, i.e. $E\tilde{\beta} = \beta$, then we only need to worry about the covariance terms. Estimates of $\beta$ obtained by ridge regression, Lasso, and LARS are all biased. This could be useful if the variance is decreased enough to give smaller MSE.

# Derived features §**3.5**

- replace $\mathbf{x}_1, \ldots \mathbf{x}_p$ with linear combinations of columns
- principal components from SVD are natural candidates
- $X = UDV$
- $z_m = Xv_m, \quad m = 1, \ldots, M < p$
- $z_m$ are orthogonal by construction
-

$$\hat{y}_{(M)}^{pcr} = \bar{y}\mathbf{1} + \sum_{m=1}^{M} \hat{\theta}_m z_m$$

-
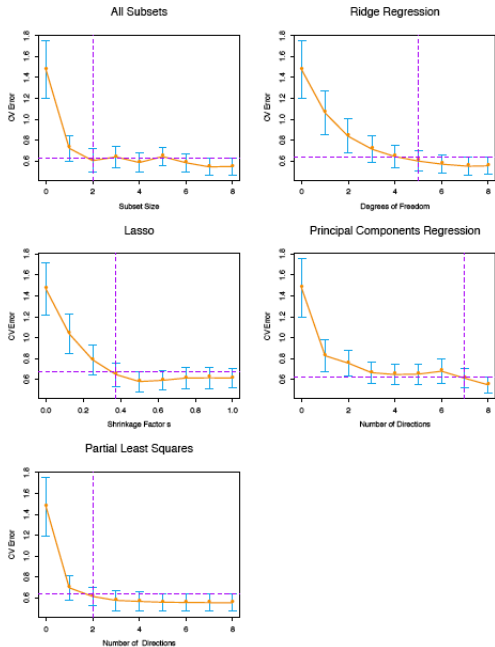
$$\hat{\theta}_m = \frac{\langle z_m, y \rangle}{\langle z_m, z_m \rangle}$$

- inputs should be scaled first (mean 0, variance 1)
- Figure 3.17

# ... derived features

- ▶ closely related method Partial least squares
- ▶ also constructs derived variables
- ▶ widely used in chemometrics, where often $p > N$
- ▶ see §3.6 for discussion

# **Thursday and next week**

- ▶ more on ridge regression and lasso in R
- ▶ construction of Table 3.3 in R test set error
- ▶ discussion of project
- ▶ HW 1 will be available
- ▶ Next week Chapter 4: §4.1 to §4.4