# Automatic Variable Selection

- Instead of choosing what goes into the model manually, we can find a program/ algorithm to do so

- Suppose we have $y$ and $q$ predictors. There are $2^q$ possible regression models

- A computer can potentially pick a "best" model, but it has to be taught what is meant by "best"

- keep in mind: bias-variance tradeoff for prediction based on learned model

  - Bias: the prediction error resulting from miss-specifying the model
  - Variance: the prediction error resulting from variations in the data used for fitting

  - Suppose we have a model $f$ to describe response $y$ from feature $z$. Given a new $z_{n+1}$ we may write

  $$y_{n+1} = \underbrace{f(z_{n+1})}_{z_{n+1}^T \beta} + \varepsilon_{n+1}$$

  - Suppose that we have past data.

this "data" $\left(y_i, z_i\right)_{i=1}^{n}$ we estimate $\hat{f}$ from

In the linear model:
$$\hat{f}(z) = z^T \hat{\beta} \quad , \quad \hat{\beta} = (Z^T Z)^{-1} Z y$$
$$z \in \mathbb{R}^p \qquad\qquad Z \in \mathbb{R}^{n \times p}$$

— We have the decomposition:
$$\mathbb{E}\left[ \left(y_{n+1} - \hat{f}(z_{n+1})\right)^2 \right]$$

$$= \mathbb{E}\left[ \left(y_{n+1} - f(z_{n+1})\right)^2 \right] \quad \text{\color{blue}{irreducible error}}$$
$$+ \left( f(z_{n+1}) - \mathbb{E}\left[ \hat{f}(z_{n+1}) \right] \right)^2 \quad \text{\color{blue}{bias}}^2$$
$$+ \mathbb{E}\left[ \left(\hat{f}(z_{n+1}) - \mathbb{E}\left[ \hat{f}(z_{n+1}) \right] \right)^2 \right] \quad \text{\color{blue}{variance}}$$

$$:= \sigma^2 + \left(\text{bias}(\hat{f})\right)^2 + \text{Var}\left( \hat{f}(z_{n+1}) \right)$$

$\uparrow$
we cannot control

Possible tradeoff based
on the way we choose $\hat{f}$

— When $f(z) = z^T \beta$ so that
$$y_{n+1} = z_{n+1}^T \beta + \varepsilon_{n+1} \qquad \mathbb{E}[\varepsilon_{n+1}] = 0, \text{ then}$$
$$\mathbb{E}\left[ \hat{f}(z_{n+1}) \right] = \mathbb{E}\left[ z_{n+1}^T \hat{\beta} \right] = z_{n+1}^T \mathbb{E}[\hat{\beta}]$$

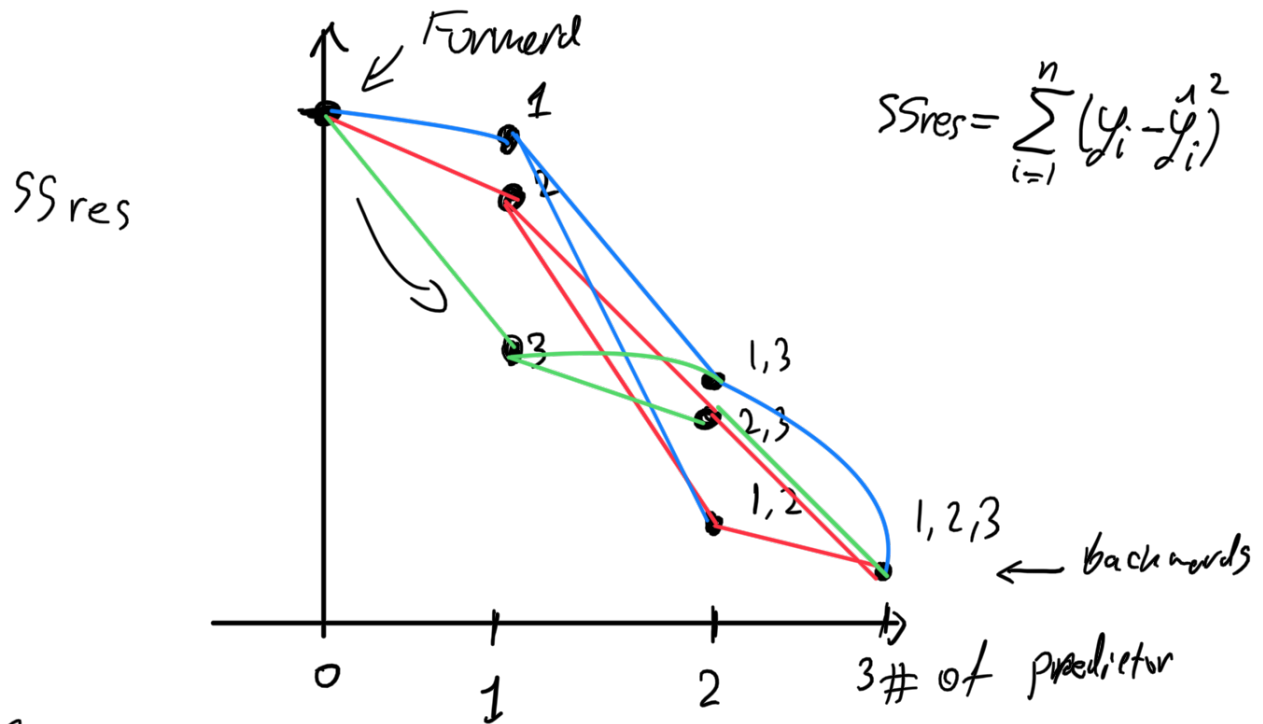$$= z_{n+1}^T \beta = \mathbb{E}[y_{n+1}] = f(z_{n+1})$$

Namely, the bias is zero.

It may not be zero anymore provided:

We don't have/use an unbiased estimator for $f$. Typically, this can happen because we don't know precisely $f$
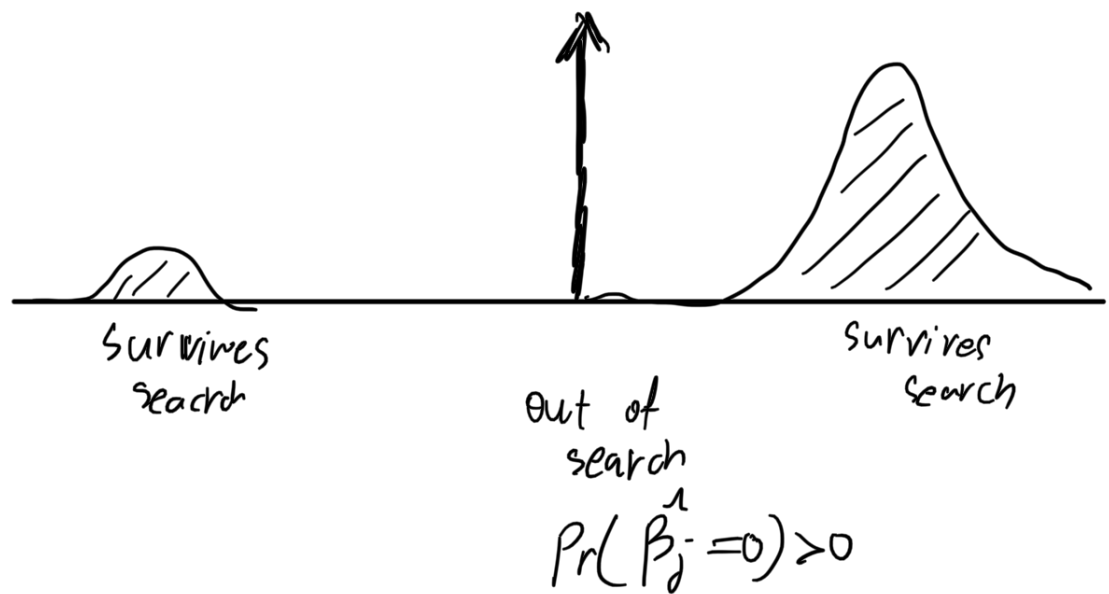
---

# All Combinations

Example: 3 predictors



SSres

$$SSres = \sum_{i=1}^{n} (y_i - \hat{y}_i)^2$$

Summary:

| σ #Pred | Best | SSres |
|---------|----------|-------|
| 0 | ∅ | 10 |
| 1 | {3} | 5 |
| 2 | {1, 2} | 2 |
| 3 | {1, 2, 3} | 1 |

- <u>Forward approach</u>: we start at ∅ and add the best predictor if the new model is statistically significant; we stop otherwise (using F-test for extra sum of squares)

- <u>Backward approach</u>: we start with all predictors, and drop the least significant / the one that leaves you with maximal $SS_{res}$ if the old model is not statistically significantly better than the new one; otherwise, stop

- Issues:
  - these are greedy approaches (a hybrid approach that "looks" into the future also makes sense)

- Both approaches usually disagree on the stoping point

- We cannot use p-values, confidence intervals, t-stats in the post selection model: "the quiet scandal"

- The distribution of $\hat{\beta_j}$ when we take the model selection process into account:



Survives search      Out of search      Survives search

$$Pr(\hat{\beta_j} = 0) > 0$$

"The quiet scandal" [Leo Breiman 92]

- Why use stepwise approaches:

  - Models can be cross-validated to measure accuracy

  - Helps leave out variables which

...reips    ...eave ... ...a.apes ...;...
Can save time & money

## Cross Validation

- split the data into $k$ smaller sets;
  leave one aside; fit a family of models
  based on the other $k-1$ sets; evaluate
  accuracy over left-out set; repeat for
  all $i=1,\dots,k$; average accuracy;
  pick the model with best averaged accuracy.

- Special case $k=n$ called leave-one-out

- Let $\hat{y}^{(i)}$ be LS predictor when
  $z_i$ is left out:

$$\hat{y}_i^{(i)} = z_i^T \hat{\beta}^{(i)}$$

where $\hat{\beta}^{(i)}$ depends on $\{(z_j, y_j)_{j \neq i}\}$
( fitted based on $Z$ with $i$-th row removed)

- Cross Validation ( CV) error:

$$CV(model) = \sum_{i=1}^{n} (y_i - \hat{y}_i^{(i)})^2$$

- Why use CV:

  - Generally, knocks out models
    that overfit
  - Parallels our goal of predicting
    new values

- Issues: can computationaly prohibitive

## CV in Linear Models

- We have a "short cut"

$$\hat{y}_i = H_{ii} y_i + (1 - H_{ii}) \hat{y}_i^{(i)}$$

$$H = Z (Z^T Z)^{-1} Z^T$$

$\left(\begin{array}{l}\text{proof is based on the Sherman-Morrison}\\ \text{formula for the inversion of matrix + rank one mat.}\end{array}\right)$

so $\quad \hat{y}_i^{(i)} = \dfrac{\hat{y}_i - H_{ii} y_i}{1 - H_{ii}}$

all left out prediction are evaluated using
one regression instead of $n$

- The residuals:

$$y_i - \hat{y}_i^{(i)} = \frac{y_i - \hat{y}_i}{1 - H_{ii}} = \frac{\hat{\varepsilon}_i}{1 - H_{ii}}$$

$$H_{ii} \leq 1$$

- Overall:

$$CV = \sum_i^n \left(y_i - \hat{y}_i^{(i)}\right)^2 = \sum_{i=1}^n \frac{\hat{\varepsilon}_i^2}{(1 - H_{ii})^2}$$

$$\boxed{\text{Penalty on number of variables}}$$

- Akaike's Information Criterion (AIC)

$$AIC := n \cdot \log\left(\frac{SS_{res}}{n}\right) + 2P$$

- Bayes Information Criterion (BIC)

$$BIC = n \log \frac{SS_{res}}{n} + p \log n$$

- We pick the model that minimizes AIC or BIC

- Considerations:

  - AIC & BIC are more "principled" approaches than forward/backward

- For $\log(n) > 2$, BIC penalizes predictors more severly.

- BIC is better at getting the "right" model correctly

identifies non-zero variables when true model is linear as $n \to \infty$

- AIC is more accurate in making predictions

- We should not use confidence intervals and t-tests after selecting models with AIC or BIC