Again on bias - variance trade-off
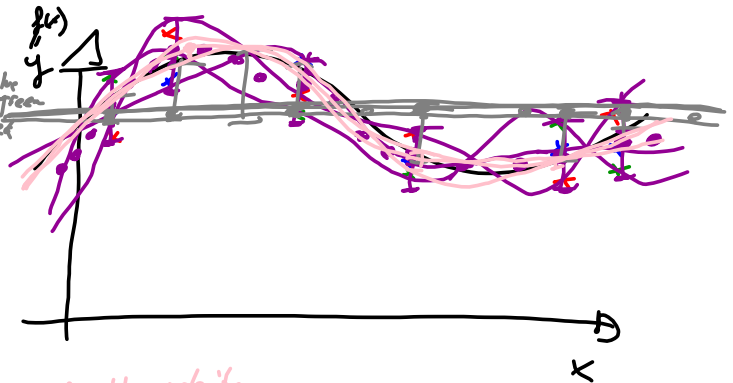- decompose squared error in bias² and variance

$$E\left[(\hat{f}(x) - f(x))^2\right]$$

$f(x)$ true value
$\hat{f}(x)$ our estimator → $\hat{y}$

$$= E\left[(\hat{f}(x) - \underline{E[\hat{f}(x)]} + \underline{E[\hat{f}(x)]} - f(x))^2\right]$$

$$= E\left[(\hat{f}(x) - E[\hat{f}(x)])^2\right] + E\left[(E[\hat{f}(x)] - f(x))^2\right] +$$

$$+ 2 E\left[(\hat{f}(x) - E[\hat{f}(x)])(E[\hat{f}(x)] - f(x))\right]$$

$$= Var(\hat{f}(x)) + (E[\hat{f}(x)] - f(x))^2 + 2(E[\hat{f}(x)] - f(x)) E[\hat{f}(x) - E[\hat{f}(x)]]$$

$$= Var(\hat{f}(x)) + bias^2 + 0$$

Empirical visualization
$$Var(\hat{f}(x)) = E[(\hat{f}(x) - E[\hat{f}(x)])]^2$$
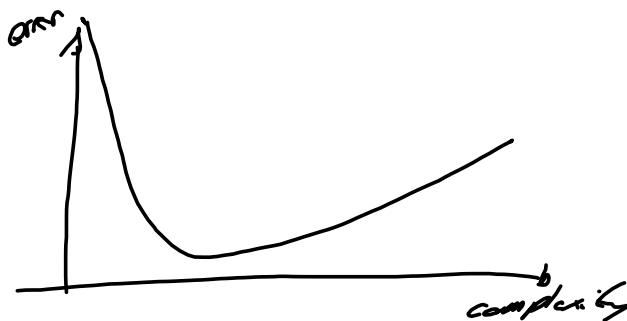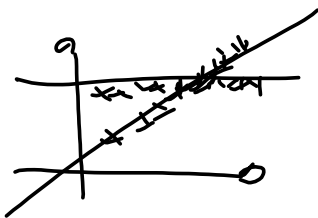$$bias^2 = (E[\hat{f}(x)] - f(x))^2$$



$\hat{f}(x)$ complex    bias very small
large variance

$\hat{f}(x)$ simple    small variance
large bias

$\hat{f}$ right complexity     best trade-off
between bias and variance

# Methods for model selection

GOAL : find the "best" model among all possible models explaining/predicting
      $y$ given $X$

Usually, we have a vector of response $y = (y_1, ..., y_n)^T$

and a design matrix
matrix of covariates
$$X = \begin{pmatrix} 1 & x_{11} & \cdots & x_{1p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \cdots & x_{np} \end{pmatrix} = (\underline{1}, X_1, ..., X_p)$$

$\longrightarrow$ $n \times (p+1)$ matrix $\Big\{$ $n$ observations
                                     $p$ variables  (+1 for the intercept)

    intercept: not strictly necessary, but
           most often used in practice

    if we use $y - \bar{y}$ as response $\longrightarrow$ intercept $= 0$

Only considering linear effects, with $p$ variables there are $2^p$ possible models to
relate $X$ and $y$

complexity $\Big\downarrow$    $\underline{y = \beta_0 + \varepsilon}$      $\hookleftarrow$ null model   (only intercept)

             $\vdots$

         $\underline{y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + ... + \beta_p x_p + \varepsilon}$    $\hookleftarrow$ full model (all variables included)

- find the best model :  1- Best Subset Selection
  - fit all the models
  - find the model that minimizes an information criterion
  - it is very computationally heavy $\hookleftarrow$ we need to fit $2^p$ models

- alternatively : simplified procedure (mostly used in practice)
  - backward elimination
  - forward selection
  - stepwise selection
  - stepback selection

## Backward elimination

→ start with the full model
→ remove one by one the least important variables (those which increase the least the residual variance)
→ proceed until the removal of the least important variable increase the chosen information criterion

$$IC = \text{deviance} + \text{penalty}$$



## Forward selection

→ start from the null model
→ add each time the variable that reduces the residual variance the most (most useful variable to explain the variability of $y$)
→ stop when adding a new variable results in an increase in the IC

$$IC = \text{deviance} + \text{penalty}$$

Stepwise selection : like forward selection, but at each step we allow the removal of a variable previously added.

Stepback selection : like backward elimination, but each step we allow the re-adding of a variable previously eliminated.

Advantages of backward elimination / stepback selection
• we start with a legitimate model
• better handling of correlation between variables

Advantages of forward selection / stepwise selection
• works also for $p > n$
• in general, much better for large $p$

# observations : $n$
# covariates : $p$

Model selection techniques:
- work in general situations (linear models, GLM, GAM, Cox model, ...)
- easy to implement, used a lot in practice
- various issues as: - underestimation of the s.e.
                     - re-use of data
                     - multiple testing

$$\hat{\beta}_B X_B + \hat{\beta}_C X_C + \hat{\beta}_D X_D + \hat{\beta}_E X_E$$

categorical variables $X \begin{cases} A \\ B \\ C \\ D \\ E \end{cases}$

$X_1 = A$
$X_2 = C$
$X_3 = D$
⋮

| intercept | | $X_C$ | $X_D$ | $X_E$ |
|---|---|---|---|---|
| $X_A$ | $X_B$ | | | |
| 1 | 0 | 0 | 0 | 0 |
| 0 | 0 | 1 | 0 | 0 |
| 0 | 0 | 0 | 1 | 0 |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋮ |

deviance $+ 2\hat{p}$

using AIC is approximately like testing at level $\alpha = 0.157$

# Principal Component Analysis / Regression

- we want to go towards a simpler (smaller $p$) model
  - → variable selection : remove useless variables;
- alternatively :
  - construct new variables as linear combinations of the original
  - try to keep as much of the original information in as **less** new variables as possible

Consider a matrix $X$ of data, suppose with mean $0$ and variance $\Sigma$

e.g. $\quad X \sim N_p(0, \Sigma)$

$\hookrightarrow$ if mean $\neq 0$, we can always center the variables $x_i = x_i - \overline{x_i}$

We want new variables $z$ on the form $\alpha X$

$$Z = \alpha X \implies Var(z) = \alpha^T \Sigma \alpha$$

The goal is to find $\alpha$ which makes $Var(z)$ the largest among all normalized linear combinations of the columns of $X$

$$\max_{\alpha} Var(z) = \max_{\alpha} \alpha^T \Sigma \alpha$$

$$\text{subject to } \|\alpha\| = 1$$

Once we did this (we obtained the $\alpha$, let us say $\alpha_1$) we look for another linear combination, orthogonal to the first one

$$\max_{\alpha} Var(z) = \max_{\alpha} \alpha^T \Sigma \alpha$$

$$\text{subject to } \|\alpha\| = 1 \quad \text{and} \quad \alpha^T \alpha_1 = 0$$

And so on...

Graphical interpretation ( in 2 dimensions, $X = (x_1, x_2)$ )



$z_1 \leftarrow$ first principal component
$z_1 = \alpha_1^{(1)} x_1 + \alpha_1^{(2)} x_2$

$z_2 \leftarrow$ second principal component

Mathematically, this can be done by using the spectral decomposition of the matrix $\Sigma$

$$\Sigma = (\alpha_1 \ldots \alpha_p) \begin{pmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_p \end{pmatrix} \begin{pmatrix} \alpha_1 \\ \vdots \\ \alpha_p \end{pmatrix} \qquad \lambda_1 \geq \lambda_2 \geq \ldots \geq \lambda_p$$

$\underset{\text{eigenvectors}}{\uparrow} \qquad\qquad\qquad \underset{\text{eigenvalues}}{\uparrow}$

$Z_1 = \alpha_1 X$

$\quad \hookleftarrow$ first vector of principal loadings

$\hookrightarrow$ first principal component

$\lambda_R$ is a measure of the original variance explained by the first principal component

- each $k$-th principal component explains $\dfrac{\lambda_k}{\sum_{k=1}^{p} \lambda_k}$ of the original variance

- the first $\underline{K}$ principal components explain $\sum_{k=1}^{K} \lambda_k \Big/ \sum_{k=1}^{p} \lambda_k$   `   `   `   `

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad \underset{\substack{\uparrow \\ \text{new variables}}}{}$

useful to select how many of the principal components to use in a model

e.g. : 90% of the original variance $\longrightarrow$ 4 principal components
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \underset{J}{\smile}$

PCR   $\quad y = \beta_0 + \beta_1 \underline{z_1} + \ldots + \beta_J \underline{z_J} + \varepsilon$

- very useful in case of large $p$ (only $k < p$ components are used)
  $\Rightarrow$ it can be used when $p > n$
- removing less useful principal components (those that do not explain anything of the original variance) reduce the model complexity
- used to solve collinearity problems (PC are orthogonal)
- used for graphical purposes (usually the first two principal components are plotted)



NOTES
- obviously, $\Sigma$ is unknown, so we need to use $\hat{\Sigma}$ in the computations;
- there are some drawbacks:
  - reduced interpretability
  - reduced portability wrt variable selection

$\underset{n \times p}{X} \xrightarrow[\text{PC}]{} Z$ $\qquad\qquad \hat{y}_{ols} = X\hat{\beta}$ $\qquad\qquad \hat{y}_{ols} = \hat{y}_{PC}$

$\qquad\qquad\qquad\qquad\qquad \hat{y}_{PC} = Z\hat{\gamma}$ $\qquad$ when we use all the principal components in our PCR $(J = p)$