# Linear regression *

## Simple regression

• Residual sum of squares $(RSS) = e_1^2 + \ldots + e_m^2 = (y_1 - \hat{B}_0 - \hat{B}_1 x_1)^2 + \ldots + (y_m - \hat{B}_0 - \hat{B}_1 x_m)^2$

• least - square minimize RSS with $\hat{B}_1 = \dfrac{\sum_{i=1}^{m}(x_i - \bar{x}_m)(y_i - \bar{y})}{\sum_{i=1}^{m}(x_i - \bar{x})^2}$    $\hat{B}_0 = \bar{y} - \hat{B}_1 \bar{x}$

• $\hat{B}_1$ and $\hat{B}_0$ are unbiased (does not systematically over or under estimate the true parameter).

→ In general sample mean unbiased $\hat{u} = \bar{y} = \frac{1}{m}\sum_{i=1}^{m} y_i$ and $SE(\hat{u})^2 = \dfrac{\sigma^2}{m}$, where $\sigma$ is the SD of each realization $y_i$ of $Y$. $SE(\hat{u})$ tells us how $\hat{u}$ differs from actual value of $u$.

$SE(\hat{B}_0)^2 = \sigma^2\left[\dfrac{1}{m} + \dfrac{\bar{x}^2}{\sum_{i=1}^{m}(x_i - \bar{x})^2}\right]$ ,   $SE(\hat{B}_1)^2 = \left[\dfrac{\sigma^2}{\sum_{i=1}^{m}(x_i - \bar{x})^2}\right]$ where $\sigma^2 = Var(\varepsilon)$

supposing common variance $\sigma^2$ and uncorrelated.

• Notice in the formula that $SE(\hat{B}_1)$ is smaller when the $x_i$ are more spread out; intuitively we have more leverage to estimate the slope when this is the case. We also see that $SE(\hat{B}_0)$ would be the same as $SE(\hat{u})$ if $\bar{x}$ were zero (in which case $\hat{B}_0 = \bar{y}$).

• $\sigma$ unknown ⇒ Residual standard error (RSE) $= \sqrt{RSS/(m-2)}$ ⇒ average amount the response will deviate from the true regression line. ⇒ measure lack of fit of the model.

• A 95% confidence interval is defined as a range of values such that with 95% probability, the range will contains the true unknown value of parameters. $\hat{B}_1 \pm 2 \cdot SE(\hat{B}_1)$

• $H_0$: there is no relationship between $X$ and $Y$ : $B_1 = 0$ ⇒ $Y = B_0 + \varepsilon$ (null)

  $H_a$: there is some " " " " " $B_1 \neq 0$ ⇒ $Y = B_0 + B_1 x + \varepsilon$ (alternative)

we need to test wether $\hat{B}_1$ is sufficiently far away from 0. How far is enough? (because if $SE(\hat{B}_1)$ is high it has to be farther than if $SE(\hat{B}_1)$ is low). so we use

t-statistic ⇒ $t = \dfrac{\hat{B}_1 - 0}{SE(\hat{B}_1)}$ (the number of SD that $\hat{B}_1$ is from 0). with $m-2$ degree freedom. (same as gaussian for $m \geqslant 30$)

we need to compute the probability of observing any number equal to $|t|$ or larger, assuming $B_1 = 0$ ⇒ this is the p value. Small p-value (< 5% or 1%) means it is unlikely to observe association between predictor and response due to chance. ⇒ reject null hypothesis

### Accuracy of the model (once rejecting the null)

RSE but also $R^2$ = the proportion of variance explained $(0 \leq R^2 \leq 1)$  $R^2 = \dfrac{TSS - RSS}{TSS} = 1 - \dfrac{RSS}{TSS}$ where $TSS = \sum(y_i - \bar{y})^2$ is total sum of squares. which is the variability inherent of the response before the regression is performed).

(TSS-RSS) amount of variability explained by the regression). $R^2$ measures the proportion of variability in Y that can be explained using X. $R^2 = 1$ good. But hard to determine the right $R^2$ = domain knowledge).

## Multiple linear regression

- $H_0: B_1 = B_2 = \cdots = B_p = 0$   $H_a$: at least one $B_g$ is non zero. This hypothesis is performed using F statistic $F = \dfrac{(TSS-RSS)/p}{RSS/(m-p-1)}$. if the linear assumptions are correct $E\{RSS/(m-p-1)\} = \theta^3$ and provided $H_0$ is true $E\{(TSS-RSS)/p\} = \theta^2$. $\Rightarrow$ iff no relationship F is close to 1, 

$\qquad\qquad$ " $> \theta^2$ iff relationship and $F > 1$

How large does F need to be? if $m$ is big small deviation is strong evidence against $H_0$. if $m$ is small we need a strong deviation. But in all case, according to $m$ and $p$, we look at the p-value of F. (because follows f-distribution).

- test a particular sublist of q of the coefficients are zero. we fit a new model without the q and get $RSS_0$ then $F = \dfrac{(RSS_0 - RSS)/q}{RSS/(m-p-1)}$

- In multiple regression, we get t-statistic and p-value for each individual (whether it is related to the response). It is equivalent of the F-statistic that omits that single variable, leaving all the others in.

- It seems that if any one of the individual p-value is small, at least one is related to the predictors, but this is flawed. $\Rightarrow$ i.e if $p = 100 \Rightarrow B_1 = \cdots B_{100} = 0$, about 5% will be below 0.05 by change. But F-statistic doesn't suffer from that as it adjusts from the number of predictors. Hence if $H_0$ is true, there is only a 5% change that the F-statistic will result in a p-value below 0.05, regardless the number of predictors. But this approach only when $p$ is small, and small relative to $m$. otherwise use forward selection or high dimensional method.

- Method in multiple regression $\longrightarrow$ got F-statistic

$\qquad\qquad\qquad\qquad\qquad\qquad\longrightarrow$ p-value big (of F) $\rightarrow$ no relationship
$\qquad\qquad\qquad\qquad\qquad\qquad\longrightarrow$ P-value of F is small

$\qquad\qquad\qquad\qquad$ if P is big $\longleftarrow\quad\searrow$ if p is small we can select individually the most promising (but not the best way).

variable selection: we can use Mallow's CP, AIC, BIC, adjusted $R^2$ (R6). But with p predictor we get $2^p$ differents models so there is easier way.

**forward selection**
begin with a null model of an intercept and no predictor. we then fit p simple linear regression and add to the null model the variable that result in the lowest RSS. we then keep going till a certain threshold.

**Backward**
we start with all variables, and remove the one with largest p value. And repeat

**Mixed**
Start with nothing, add till p-value for one of the variables ↗, then we remove it. we repeat till small p-value for all.

warning: $R^2$ always increase (closer to 1) as we add variables, because always reduce the residual sum of squares on the training data. Also $RSE = \sqrt{\dfrac{1}{m-p-1} RSS}$ then we can get ↗ RSE if the decrease in RSS is small relative to the increase in p.