# Statistical Regression

→ Regression is a means of exploring the **variation** in some quantity

EXPLAINED    UNEXPLAINED

→ E.g.   ICE CREAM SALES

- Temperature
- Rainfall
- School holidays
- ?

Regression will be able to quantify How much can be **explained** & how much is **unexplained**

---

→ Why machine learning can be treated as a statistical inference problem?

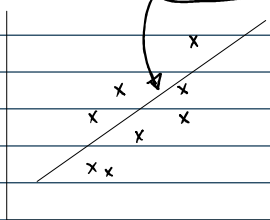e.g.  lets try to explain the variation in ice cream sales using only one variable (daily temperature)

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$ ⟶ Population regression equation

→ Regression aims to:

① Estimating the $\beta$'s
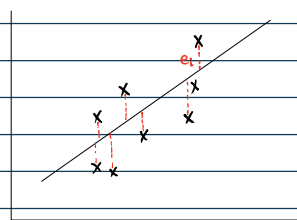
② Quantifying the errors

**Question:** Do we have all the data there is about daily temperatures and ice cream sales? **NO**, we will mostly be working with sample data, so we can never really know the true values of $\beta_0$ & $\beta_1$, we can only estimate them based on sample data.

---

→ Sample regression line is given by: $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_1$  ( this describes the <u>line itself</u> )
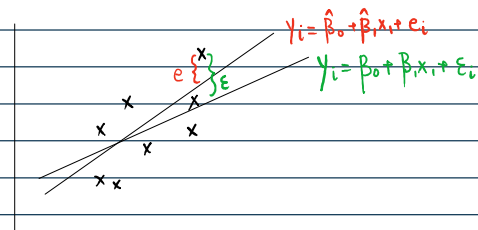
→ Each sample will yield a different sample regression line.

---

→ Value of y for each observation is given by : $y_i = \hat{\beta}_0 + \hat{\beta}_1 x_1 + e_i$

---

→   $\underbrace{y_i = \hat{\beta}_0 + \hat{\beta}_1 x_1 + e_i}_{f'(x)}$   IS AN ESTIMATE OF   $\underbrace{y_i = \beta_0 + \beta_1 x_1 + \varepsilon_i}_{f(x)}$

↳ Estimated relationship              ↳ True relationship

$y_i = \hat{\beta}_0 + \hat{\beta}_1 x_1 + e_i$
$y_i = \beta_0 + \beta_1 x_1 + \varepsilon_i$

$y = f'(x) + e_i$  (estimated relationship)

reducable        irreducable
$f(x) - f'(x)$

---

Sum of Squares due to Regression (SSR) (Explained)
$\sum (\hat{y}_i - \bar{y})^2$

→ Regression begins to look a lot like ANOVA i.e. the total sum of squares is partitioned b/w SSE & SSR.

→ Total variance (SST)
$\sum (y_i - \bar{y})^2$

Sum of squares due to Error (SSE) (still unexplained)
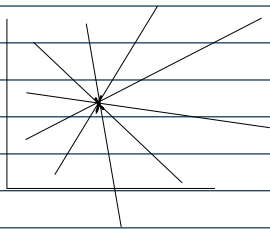$\sum (y_i - \hat{y})^2$
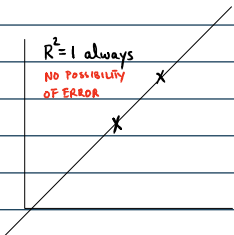
$$SST = SSR + SSE$$

→ coefficient of determination $(R^2) = \dfrac{SSR}{SST}$ ( proportion of the variation in dependent variable explained by the independent variable)

$(-\infty, 1]$

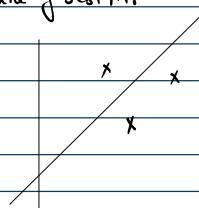→ How many observations required to perform regression

one?



two?

$R^2 = 1$ always
NO POSSIBILITY OF ERROR



→ Remember equation of regression, $y = \beta_0 + \beta_1 X_1 + \varepsilon$.
Regression requires a possibility of error.

→ With 3 observations error is introduced, NOW we can find the line of best fit.
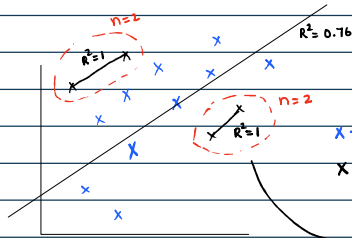


$n = 3$
$df = 1 \; (n - (k+1))$

↑
# of observations

→ $(k+1)$ represents the number of parameters being estimated (including the intercept)

→ # of independent variable

→ With $R^2 = 1$ always, the strength of the relationship b/w x & y can't be assessed. we are not looking at the big picture.

$$df\_total = df\_model + df\_residual$$
$$n - 1 = k + n - (k+1)$$

degree of freedom for the residuals



$R^2 = 0.76$

x – The big picture (can be used to make predictions in the future) ✓

x – narrow vision (memorized data instead of generalized relationship) ✗

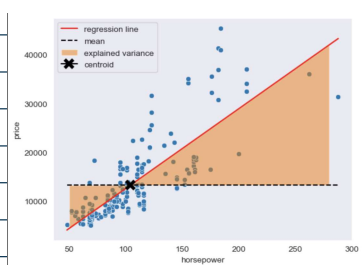→ No point in making predictions about data that can't move/vary.

→ By adding more independent variables, degrees of freedom are reduced.
Possibility of error in the model is reduced, i.e. $R^2$ continues to increase as we add more variables, we a fooling ourselves that the model is better, when in reality it's not. This is resolved by using the Adjusted $R^2$ which takes into account the reduced degrees of freedom.

$$Adj\, R^2 = 1 - \left[\dfrac{(1-R^2)(n-1)}{n-k-1}\right]$$

→ 1) we are able to explain a lot of variation in the dependent variable
When the data is free then we have a good model. More degrees of freedom (larger sample size) gives us a greater statistical power.

Regression Analysis

→ e.g. Car Price VS horsepower



**Linear Regression ▾**

Pearson correlation

Model Summary – price

| Model | R | $R^2$ | Adjusted $R^2$ | RMSE |
|---|---|---|---|---|
| $H_1$ | 0.808 | 0.651 | 0.651 | 4716.946 |

Note. Null model includes horsepower

66.3% of Total variance explained

**ANOVA**

$k$ (# of independent variables) (df model)

| Model | | Sum of Squares | df | Mean Square | F | p |
|---|---|---|---|---|---|---|
| $H_1$ | Regression | 8.503×10⁹ | 1 | 8.503×10⁹ | 382.161 | < .001 |
| | Residual | 4.517×10⁹ | 203 | 2.225×10⁷ | | |
| | Total | 1.302×10¹⁰ | 204 | | | |

SSR — SSE — MSE — SST

$\dfrac{MSR}{MSE}$

$F_{k, df\_res, df\_MSE}$ right tail test

Note. Null model includes horsepower

# of observations − 1 (n−k−1)
# of observations (n)

**Coefficients ▾**

(RMSE)
→ The standard error is the standard deviation of the error term, $\varepsilon$.

It is the average distance an observation falls from the regression line in units of the dependent variable.

We can think of s as a measure of how well the regression model makes prediction.

One of the assumptions for regression states that for a given x the error terms are normally distributed with $\mu = 0$ & $\sigma = RMSE$. Can also be used for outlier detection

$$RMSE = S = \sqrt{MSE}$$

At the top, a statistics output with annotations:

| | | SS | | | | |
|---|---|---|---|---|---|---|
| $H_0$ | Regression | **SSR** $8.503 \times 10^9$ | ① | $8.503 \times 10^9$ | $382.163$ | $< .001$ |
| | Residual | **SSE** $4.517 \times 10^7$ | 203 | $2.225 \times 10^7$ | **MSE** | |
| | Total | **SST** $1.302 \times 10^{10}$ | 204 | | | |

Note. Null model includes horsepower

from $\hat{\beta}_0 \hat{\beta}_1$, $\alpha$ b/w observations -1 b/w (n-k-1)
observations -1 observations (n)
we only require (n-1)

Coefficients ▼

$F_c$, $df_{num}$, $df_{MSE}$
right tail test

$\frac{E_c}{MSE}$

| | | | | | | | 95% CI | |
|---|---|---|---|---|---|---|---|---|
| Model | | Coef | Standard Error | | t | p | Lower | Upper |
| $H_0$ | (Intercept) | $-3721.761$ | $929.849$ | | $-4.003$ | $<.001$ | $-5555.163$ | $-1888.360$ |
| | horsepower | $163.263$ | $8.351$ | | $19.549$ | $<.001$ | $46.796$ | $179.731$ |

$t = \frac{coef}{std\ error}$

$t_{\alpha/2}$, $df_{MSE}$

$t_{1-\alpha/2}$, $df_{MSE}$

$H_0: \beta_k = 0$

$H_a: \beta_k \neq 0$

two tail test

$\dfrac{RMSE}{\sqrt{\sum(x_i - \bar{x})^2}}$

$\beta_k \pm t_{\alpha/2}\, std\ error$

model makes prediction.

$$RMSE = S = \sqrt{MSE}$$

Does a statistically significant linear relationship exist b/w the independent & dependent variables? Is the overall F-test or t-test (in simple linear regression these are actually the same thing) significant. $\boxed{F = t^2}$ only for SLR.

$H_0: \beta_1, \beta_2, \beta_3 ... \beta_k = 0$ Expect Intercept

$H_a:$ At least one regression coefficient is not equal to 0

MSE ($s^2$) is an estimate of $\sigma^2$ the variance of the error, $\varepsilon$. In other words, how spread out the data points are from the regression line.

$$s^2 = MSE = \frac{SSE}{df_{SSE}}$$

→ Would a regression analysis offer anything more than the $\bar{y}$ model? Using this nonregression model ($\bar{y}$) as a worst case, we can analyse the regression line to determine whether it adds a more significant amount of predictability of y than the $\bar{y}$ model.

→ if the slope is not different from 0, the regression line is doing nothing more than the $\bar{y}$ line in predicting y.

→ Remember that for each sample we will obtain a different slope value. The Question if all the pairs of data points for the population were available, would the slope of the regression line be different from 0?

→ For simple linear regression $\boxed{F = t^2}$

Note: We can only rely on the numbers above IF CERTAIN ASSUMPTIONS HOLD

## Prediction Intervals for Linear Regression

$$\hat{y} \pm t_{\frac{\alpha}{2}} S_e \sqrt{1 + \frac{1}{n} + \frac{n(x_0 - \bar{x})^2}{n(\sum x_i^2) - (\sum x_i)^2}}$$

| Symbol | What it Means | How To Find It |
|---|---|---|
| $\hat{y}$ | Predicted y | Plug given x into LSRL equation |
| $t_{\frac{\alpha}{2}}$ | Critical Value | $T.INV.2T(alpha, df)$ where $df = n - 2$ |
| $y - \hat{y}$ | Residual | Subtract the predicted value for each x (using LSRL) from the actual value for each x |
| SSE | Sum of Squared Errors | The sum of the squares of the residuals, $\sum(y_i - \hat{y}_i)^2$ |
| $S_e$ | Standard Error of our Sample | The root of the SSE divided by the degrees of freedom, $\sqrt{SSE/df}$ |
| $n$ | Number of Pairs | The count of the number of pairs of data |
| $\bar{x}$ | Mean of x | The sum of x-values divided by the number of x-values, $\sum x_i / n$ |
| $\sum x_i$ | Sum of x's | The sum of all given values of x |
| $\sum x_i^2$ | Sum of squares of x's | The sum of the squares of each x-value |