# Simple Linear Regression

**What Is Regression?**
- Fitting a statistical model to data, in particular linear models.
- Understand the errors in the data and how they affect the estimated model.

**Why Do Regression?**
- Describe the relationship between two variables.
- Predict the value of one variable given the value of another.
- Learn how to control future values of a variable by controlling values of variables it's related to.

**Regression Model**

$$Y = \beta_0 + \beta_1 X + \varepsilon$$

For particular points, $Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i, i = 1, \ldots, n$.
We assume $X$ is not random; if it is, condition on it.
If we have an infinite number of observations $Y$'s for each value $X$, the mean of the $Y$'s is $\bar{Y} = E[\bar{X}]$; if $X$ is random, then $E[Y|X=x]$.

**Fitted Value**
$E[Y] = \beta_0 + \beta_1 X$.
For each $X_i$, the fitted value $\hat{Y}_i = b_0 + b_1 X_i$. Deviation from the line $Y_i - \hat{Y}_i = e_i$ is the residual.

**Gauss-Markov Conditions**
1. $E(\varepsilon_i) = 0 \Rightarrow E(Y_i) = \beta_0 + \beta_1 X_i$.
2. $var(\varepsilon_i) = \sigma^2$ ( $E(\varepsilon_i^2) = \sigma^2$ ).
3. $\varepsilon_i$'s are uncorrelated ( $E(\varepsilon_i \varepsilon_j) = 0$ ).

Or equivalently,
1. $E(Y_i) = \beta_0 + \beta_1 X_i$.
2. $var(Y_i) = \sigma^2$.
3. $Y_i$'s are uncorrelated.

## METHOD OF LEAST-SQUARES

Define $S = \sum_{i=1}^{n} (y_i - b_0 - b_1 x_i)^2$ to be the sum of square residuals (vertical deviations from the line).

Since $\beta_0$, $\beta_1$, $\sigma^2$ are parameters (constants we usually don't know), they are estimated by:
- $b_0 = \bar{y} - b_1 \bar{x}$,

- $b_1 = \dfrac{\sum_{i=1}^{n} x_i y_i - n\bar{x}\bar{y}}{\sum_{i=1}^{n} x_i^2 - n\bar{x}^2} = \dfrac{\sum_{i=1}^{n} (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^{n} (x_i - \bar{x})^2}$,

- $s^2 = \sum_{i=1}^{n} \dfrac{e_i^2}{n-2}$ .

**Properties of Least Square Estimators**

1. Unbiased: $E(b_0)=\beta_0$ and $E(b_1)=\beta_1$ .
2. Linear in $Y_i$ 's: can find constants $c_i$ , $d_i$ such that $b_0=\sum c_i Y_i$ and $b_1=\sum d_i Y_i$ .

**Gauss-Markov Theorem**
Least squares estimate are the best (minimum variance) linear unbiased estimators (BLUE).

**Properties of Fitted Regression Line**

1. $\sum_{i=1}^{n} e_i=0$ , where $e_i=Y_i-\hat{Y}_i$ are the residuals.

2. $\sum_{i=1}^{n} e_i^2$ is a minimum.

3. $\sum_{i=1}^{n} e_i X_i=0$ .

4. $\sum_{i=1}^{n} \hat{Y}_i=n\bar{Y}$ , that is the average of the predicted values ( $\hat{Y}_i$ ) is the same as the average of the data values of $Y$.

5. $\sum_{i=1}^{n} e_i \hat{Y}_i=0$ .

Note: $\sum_{i=1}^{n} e_i Y_i \neq 0$ .

## CORRELATION
Correlation gives rv's $X$ and $Y$ a symmetric measure of the strength of their linear relationship.

**Pearson Product-Moment Correlation**

$$r=\frac{\sum_{i=1}^{n}(X_i-\bar{X})(Y_i-\bar{Y})}{\sqrt{\sum_{i=1}^{n}(X_i-\bar{X})^2-\sum_{i=1}^{n}(Y_i-\bar{Y})^2}}$$

This is which is the MLE of $\rho=\dfrac{cov(X,Y)}{\sqrt{var(X)var(Y)}}$ when $X$ and $Y$ have a bivariate normal distribution.

Note: $b_1=r\dfrac{s_y}{s_x}$ , where $s_x$ and $s_y$ are the standard deviations for $X$ and $Y$.

## REGRESSION ANALYSIS OF VARIANCE

**Decomposition of Sum of Squares**

$$\sum_{i=1}^{n}(Y_i-\bar{Y})^2=b_1^2\sum_{i=1}^{n}(X_i-\bar{X})^2+\sum_{i=1}^{n}e_i^2$$

- $\sum_{i=1}^{n}(Y_i-\bar{Y})^2$ = Total Sum of Squares (SSTO).
- $b_1^2\sum_{i=1}^{n}(X_i-\bar{X})^2$ = Regression Sum of Squares (SSR) – amount of variation in the $Y$'s that is explained by the least-square line.
- $\sum_{i=1}^{n}e_i^2$ = Error Sum of Squares (SSE).

**ANOVA Table**

| Source | SS | df | MS |
|---|---|---|---|
| Regression Line | SSR = $b_1^2 \sum_{i=1}^{n} (X_i - \bar{X})^2$ | 1 | $\dfrac{\text{SSR}}{1} = \text{MSR}$ |
| Residuals | SSE = $\sum_{i=1}^{n} e_i^2$ | $n-2$ | $\dfrac{\text{SSR}}{n-2} = \text{MSE} = s^2$ |
| Total | SSTO = $\sum_{i=1}^{n} (Y_i - \bar{Y})^2$ | $n-1$ | |

**Coefficient of Variation**

The coefficient of variation is $R^2 = \dfrac{\text{SSR}}{\text{SSTO}} = \dfrac{b^2 \sum (X_i - \bar{X})^2}{\sum (Y_i - \bar{Y})^2} = 1 - \dfrac{\sum e_i^2}{\sum (Y_i - \bar{Y})^2}, \ 0 \le R^2 \le 1$ . It gives the percentage of

variation in the $Y$'s explained by the regression line.

Note: $R^2 = \dfrac{\sum (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum (X_i - \bar{X})^2 \sum (Y_i - \bar{Y})^2}}$ which is the square of sample correlation coefficient between $X$ and $Y$.

# STATISTICAL PROPERTIES OF $b_0$ AND $b_1$

$b_0$ and $b_1$ are rv's because they are a function of the $Y$'s. So they have expected values and variances.

1. $E(b_0) = \beta_0$ , $\text{var}(b_0) = \sigma^2 \left( \dfrac{1}{n} + \dfrac{\bar{X}^2}{\sum (X_i - \bar{X})^2} \right)$ .

2. $E(b_1) = \beta_1$ , $\text{var}(b_1) = \dfrac{\sigma^2}{\sum (X_i - \bar{X})^2}$ .

# HYPOTHESIS TESTING

**Steps**
1. Specify the null ( $H_0$ ) and the alternative ( $H_A$ ) hypothesis for parameter of interest.
2. Calculate a test statistic whose distribution is known.
3. Estimate a p-value. Assuming $H_0$ is correct, calculate the probability of the value of the test statistic for values more extreme (belongs to $H_A$ ). P-value gives evidence against $H_0$ .
4. If p-value is small, then either $H_0$ is correct and data are just rare, or $H_0$ is incorrect. The smaller the p-value, the stronger the evidence against $H_0$ .
   $p > 0.1$ : no evidence against $H_0$ .
   $0.05 < p < 0.1$ : some weak evidence against $H_0$ .
   $0.01 < p < 0.05$ : moderate evidence against $H_0$ .
   $p < 0.01$ : strong evidence against $H_0$ .

**t-Test**
If $X_1, \ldots, X_n$ are iid $N(\mu, \sigma^2)$ rv's, then:

1. $\bar{X} \sim N\left( \mu, \dfrac{\sigma^2}{n} \right)$ .

2. $\bar{X}$ is an estimator of $\mu$ , $s^2 = \dfrac{\sum (X_i - \bar{X})^2}{n-1}$ is an estimator of $\sigma^2$ .

3. If we standardize $\bar{X}$ , we get $\dfrac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}} \sim N(0,1)$ . If we estimate $\sigma$ by $s$, then $\dfrac{\bar{X} - \mu}{\frac{s}{\sqrt{n}}} \sim t_n - 1$ .

**Hypothesis Test for $\beta_1$**

$b_1 \sim N\left(\beta_1, \dfrac{\sigma^2}{\sum (X_i - \bar{X})^2}\right)$ because it is a linear combination of normal rv's.

Estimate $var(b_1)$ by $\dfrac{s^2}{\sum (X_i - \bar{X})^2}$ which has $n-2$ df. Thus, $\dfrac{b_1 - \beta_1}{\sqrt{\dfrac{s^2}{\sum (X_i - \bar{X})^2}}} \sim t_{n-2}$ .

To test $H_0: \beta_1 = 0$ vs $H_A: \beta_1 \neq 0$ , use test statistic $t_{obs} = \dfrac{b_1}{\sqrt{s^2 / \sum (X_i - \bar{X})^2}}$ which is from $t_{n-2}$ if $H_0$ is true.

## CONFIDENCE INTERVALS

Let $X_1, \ldots, X_n$ be i.i.d. $N(\mu, \sigma^2)$ . We want a confidence interval for $\mu$ .

Since $\dfrac{\bar{X} - \mu}{S/\sqrt{n}} \sim t_{n-1}$ , we find $t_{n-1, \alpha/2}$ such that $P\left(\left|\dfrac{\bar{X} - \mu}{S/\sqrt{n}}\right| \leq t_{n-1, \alpha/2}\right) = 1 - \alpha$ . Then the interval $\bar{X} \pm t_{n-1, \alpha/2} \dfrac{S}{\sqrt{n}}$ is a

$100(1-\alpha)$ % C.I. for $\mu$ .

Not true: The probability that $\mu$ is in this interval is $(1-\alpha)$ .

True: C.I. calculated from repeated sample from a $N(\mu, \sigma^2)$ of size $n$ by this method will include $\mu$ $100(1-\alpha)$ % of the time.

Note: Quote C.I. to give an idea of the precision of an estimate.

## Prediction

Least square line: $\hat{Y} = b_0 + b_1 X$ .

For $X_h$ (a given value of $X$), a $100(1-\alpha)$ % C.I. for the corresponding point on the line (estimate of $E(Y_h)$ ) is

$\hat{Y}_h \pm t_{n-2, \alpha/2} s \sqrt{\dfrac{1}{n} + \dfrac{(X_h - \bar{X})^2}{S_{XX}}}$ .

Note: If the line is a good fit, $s$ is small and narrow C.I.

## Prediction Interval

We use the line to predict a new $Y_h$ given a $X_h$ .

A $100(1-\alpha)$ % prediction interval for $Y_h$ corresponding to $X_h$ is $\hat{Y}_h \pm t_{n-2, \alpha/2} s \sqrt{1 + \dfrac{1}{n} + \dfrac{(X_h - \bar{X})^2}{S_{XX}}}$

Note: This is not a C.I. since C.I. are for parameters.

## Confidence Band: Workman-Hotelling Procedure

Want to put a C.I. on the entire line, i.e. a confidence band around the line.

Replace t-critical values $t_{n-2, \alpha/2}$ in C.I. by $\sqrt{2 F_{2, n-2, 1-\alpha}}$ , where $F_{2, n-2, 1-\alpha}$ is the $100(1-\alpha)$ percentile from the $F$ distribution with 2 and $n-2$ df.

Thus, the bounds at $X_h$ is $\hat{Y}_h \pm W \operatorname{var}(\hat{Y}_h) = \hat{Y}_h \pm W \sqrt{s^2 \left( \dfrac{1}{n} + \dfrac{(X_h - \bar{X})^2}{\sum (X_i - \bar{X})^2} \right)}$ . To obtain the confidence band, join up all the C.I.'s.

## The ANOVA F-Test

If $\beta_1 = 0$ , then $E(\text{MSE}) = E(\text{MSR}) = \sigma^2$ .

Test $H_0 : \beta_1 = 0$ vs $H_A : \beta_1 \neq 0$ , use test statistic $F_{obs} = \dfrac{\text{MSR}}{\text{MSE}}$ . Since $\beta_1 \neq 0$ gives larger values of $F_{obs}$ , deviations from $\beta_1 = 0$ are only in the right tail of $F$ distribution.

# Diagnostic and Remedial Measures

## So Far

- Model: $Y_i = \beta_0 + \beta_1 X_i + \epsilon_i, n = 1, \dots, n$ .
- Assumption: Linear relationship is appropriate.
- LS line: $\hat{Y}_i = b_0 + b_1 X_i$ .
- Gauss-Markov assumptions: $E(\epsilon_i) = 0$ , $\operatorname{var}(\epsilon_i) = \sigma^2$ , $\epsilon_i$ and $\epsilon_j$ uncorrelated; can show $b_0$ and $b_1$ unbiased and BLUE.
- For testing, confidence intervals, assume $\epsilon_i$ normally distributed.

## Residual Plots

$e_i = Y_i - \hat{Y}_i = Y_i - b_0 - b_1 X_i$ estimate of $\epsilon_i$ . What do we know about the $e_i$ 's?

- $\sum e_i = 0$ , $\sum e_i X_i = 0$ , $\sum e_i \hat{Y}_i = 0$ .
- $e_i$ 's are random variables: $E(e_i) = 0$ , $\operatorname{var}(e_i) = \text{ugly}$ , $e_i$ 's not uncorrelated (since $\sum e_i = 0$ ).
- Distribution of the $e_i$ 's is normal because $e_i$ 's linear function of $Y_i$ 's.

Plots of $e_i$ 's are very useful for checking the assumptions of the model. Look for:

- Evidence that straight line is an appropriate model – no curvature, outlier, influential point.
- Constant variance.
- Normally distributed errors – otherwise p-values are wrong because test statistic doesn't have $t$ or $F$ distribution under $H_0 : \beta_1 = 0$ , and CI's have coverage other than $100(1 - \alpha)$ %; especially important for prediction intervals.
- Lack of independence (if possible to check).
- Missing predictor variables (if possible to check).

Recommended plots: $e_i$ vs $\hat{Y}_i$ , $e_i$ vs $X_i$ , normal quantile-quantile plot.

## Plot 1: $e_i$ vs $\hat{Y}_i$

Look for:

1. Evidence of not constant variance.
2. Outliers – more exaggerated on residual plots than scatter plots.
3. Evidence of curvature – more exaggerated on residual plots than scatter plots.
4. A "good" plot has just random scatter about 0.

**Plot 2:** $e_i$ **vs** $X_i$
- When model has only one predictor variable (simple linear regression), plots 1 and 2 are essentially the same because $Y_i = b_0 + b_1 X_i$ .
- As a general rule, curvature is assessed from this plot, where plot 1 is used for outliers and constant variance.

### Plot 3: Normal Quantile-Quantile Plot of Residuals
- Used for assessing normality.
- Check after the other two.

### Other Plots
- Univariate analysis of residuals (ex: box plot, stem-i-leaf plot): look for departures from normality.
- $|e_i|$ vs $Y_i$ : exaggerated non-constant variance.
- Residuals vs time or other sequence of observations.: a pattern indicates residuals may be correlated.
- Residual vs other predictor variables: want random scatter, otherwise this predictor variable should be considered in model.

## NORMAL QUANTILE PLOTS
- To assess assumption of normality of errors.
- More sensitive to departures from normality than a histogram.

### Idea of Construction
1. Order the $e_i$ 's from smallest to largest: $e_{(1)} \leq e_{(2)} \leq \cdots \leq e_{(n)}$ .
2. Get values of corresponding quantile from a stander normal distribution: $E(Z_{(1)}), \ldots, E(Z_{(n)})$ with $E(Z_{(i)}) = \Phi^{-1}\left(\frac{i}{n}\right)$ .
3. Plot $e_i$ vs $E(Z_{(i)})$ . If the normal assumption for $\varepsilon_i$ is valid, the plot should be approximately a straight line.

### Problem
$$E(Z_{(n)}) = \Phi^{-1}\left(\frac{n}{n}\right) = \Phi^{-1}(1) = \infty \ .$$

One possible adjustment is $E(Z_{(i)}) = \Phi^{-1}\left(\frac{i}{n+1}\right)$ .

Another possible adjustment (Blom) is $E(Z_{(i)}) = \Phi^{-1}\left(\frac{i-3/8}{n+3/4}\right)$ .

### Is Lack of Normality a Problem?
- Yes for prediction intervals.
- Less so for other tests and CI's.
- If tails are a little heavy but the distribution is of residuals is still symmetric, the $F$ and $t$ tests are still approximately correct.

## REMEDIAL MEASURES
What to do if assumptions are violated?
1. Abandon SLR for something else (usually more complicated).
2. For outliers, remove if shown to be mistakes, otherwise proceed with caution.
3. For influential points, fit regression on smaller range of values and report this.

4. Transform the data for non-linearity, non-normal residuals, non-constant variance.
   - If relationship is non-linear but variability is approximately constant over $X$, try to find a transformation of $X$ that makes relationship appear linear.
   - If relationship is non-linear and variability is non-constant, transform $Y$. Another possible solution is to use weighted least squares.
   - Sometimes, transformation of both variables is needed.

**Transformations to Stabilize the Variance**

If $Y$ has a distribution with mean $\mu_Y$ and variance $\sigma_Y^2$, then the mean and variance of $Z = f(Y)$ are approximately

$$E(Z) \approx f(\mu_Y) \quad \text{and} \quad \mathrm{var}(Z) \approx \sigma_Y^2 [f'(\mu_Y)]^2 \quad .$$

Now suppose $Y$ has mean $\mu$ and variance proportional to a function of $\mu$, i.e. $\mathrm{var}(Y) = \sigma_2 V(\mu) = \sigma_Y^2$ (in regression $E(Y) = \beta_0 + \beta_1 X$ ). We need a transformation $Z = f(Y)$ so that $\mathrm{var}(Z)$ is approximately constant. We need

$$[f'(\mu)]^2 = c \frac{1}{V(\mu)} \Rightarrow f'(\mu) \propto V^{-1/2}(\mu) \quad , \text{ so take } f(\mu) \propto \int V^{-1/2}(\mu) d\mu .$$

Some popular transformations:
   - If $\mathrm{var}(Y_i)$ is linear in $E(Y_i)$, then take $Z_i = \sqrt{Y_i}$ .
   - If $\mathrm{var}(Y_i)$ increases faster than linear in $E(Y_i)$, then take $Z_i = \ln Y_i$ .
   - $Z_i = \dfrac{1}{Y_i}$ is even more extreme.

**Semi-Studentized Residual Plots**

A semi-studentized residual is $e_i^* = \dfrac{e_i - \bar{e}}{\sqrt{\mathrm{MSE}}}$ . It is not "studentized" because MSE is an estimate of $\mathrm{SD}(\varepsilon_i)$, not $\mathrm{SD}(e_i)$ .

Plot these residuals on standard scale; helpful for evaluating outliers (more than 3 s.d. from 0).

# Simultaneous Inference and Other Topics

## JOINT ESTIMATE OF B₀ AND B₁

We want joint $100(1-\alpha)$ % CI for $\beta_0$ and $\beta_1$ . The probability that both CI capture their true value is less than $100(1-\alpha)$ %. Want to be $100(1-\alpha)$ % confident that both CI's capture what they are supposed to capture.

A solution is to use Bonferroni corrected CI's. In general, for joint/simultaneous CI for any $k$ parameters, the overall confidence level will be at least $100(1-\alpha)$ % if the confidence level for each parameter is $100\left(1 - \dfrac{\alpha}{k}\right)$ %

## SIMULTANEOUS ESTIMATE OF MEAN RESPONSE

We want joint CI for $E(Y)$ at more than one $X$ .
There are two choices: Bonferroni or Working-Hotelling.
   - Both of these procedures are conservative, i.e. the overall confidence level will actually be higher than $100(1-\alpha)$ %.
   - For a small number of $X$'s, Bonferroni is generally less conservative. For more, Working-Hotelling is less conservative.

**Note**

The Bonferroni procedure is also helpful for hypothesis testing when carrying out multiple tests.

For example, if you want to declare the result statistically significant if p-values are less then $\alpha$, then you will mistakenly reject $H_0$ $100\alpha$ % of the time (type 1 error). So if $\alpha = 0.05$ and doing 20 tests, on average one test will be reject even if $H_0$ is true all the time.

The Bonferroni Correction: When doing $k$ tests, reject $H_0$ only if p-value $\leq \dfrac{\alpha}{k}$. The probability of one or more type 1 error is then at most $\alpha$. Note that this is not useful if $k$ is large.

# INVERSE PREDICTION (CALIBRATION PROBLEM)

Suppose we observe $Y_0$. We want to estimate $X_0$.

### Inverse Estimator

The estimate $\hat{X}_0 = \dfrac{Y_0 - b_0}{b_1}$ is the inverse estimator.

### Inverse Prediction Interval

How good is the inverse estimator $\hat{X}_0$?

The PI formula is $\hat{Y}_h \pm t_{n-2, 1-\frac{\alpha}{2}} s \sqrt{1 + \dfrac{1}{n} + \dfrac{(X_h - \bar{X})^2}{S_{xx}}}$. Solving for $X_h$ is a mess. If we approximate $\dfrac{t^2_{n-2, 1-\alpha/2} s^2}{b_1^2 S_{xx}}$ by 0

when it is small (say < 0.05), then the approximate inverse prediction interval for $X_0$ when $Y = Y_0$ is

$\hat{X}_0 \pm \dfrac{t_{n-2, 1-\frac{\alpha}{2}}}{|b_1|} s \sqrt{1 + \dfrac{1}{n} + \dfrac{(\hat{X}_0 - \bar{X})^2}{S_{xx}}}$.

Note that $\dfrac{t^2_{n-2, 1-\alpha/2} s^2}{b_1^2 S_{xx}}$ being small means the test of $H_0 : \beta_1 = 0$ will be statistically significant.

# WHAT IF PREDICTOR X IS RANDOM?

- Can analyze the correlation between $X$ and $Y$. There exists tests and CI's for $\rho$ under the assumption that $X$ and $Y$ have bivariate normal distribution(SAS gives p-values for 2-sided test of $H_0 : \rho = 0$ with "proc corr").
- If we want to do regression, there is a clear choice for the independent and dependent variables.
- If $(X_i, Y_i)$ independent pairs from bivariate normal distribution, then the distribution of $Y_i | X_i$ is normal with mean

  $\alpha_0 + \alpha_1 X_i$ where $\alpha_0 = \mu_Y - \mu_X \rho_{XY} \dfrac{\sigma_Y}{\sigma_X}$ and $\alpha_1 = \rho_{XY} \dfrac{\sigma_Y}{\sigma_X}$ (still a straight line model), and constant variance

  $\sigma_Y^2 (1 - \rho_{XY}^2)$. So regression techniques (including inference) up to now still hold.
- If $(X_i, Y_i)$ are not bivariate normal, as long as the $X_i$'s are independent rv's whose probability distributions do not involve any of the regression parameters ($\beta_0$, $\beta_1$, $\sigma^2$) and the conditional distribution $Y_i | X_i$ is normal, then everything so far still hold, except need to talk about the distribution of $Y_i | X_i$ (ex: Gauss-Markov Condition: $E(Y_i | X_i) = \beta_0 + \beta_1 X_i$, $\text{var}(Y_i | X_i) = \sigma^2$, distribution of $Y_i | X_i$ independent).

### Measurement Error

Suppose $(X_i, Y_i)$ bivariate normal, but the $X_i$'s are observed with measurement error. Note that measurement error in the

$Y_i$ 's are absorbed in $\varepsilon_i$ 's.

We want the model $Y_i + \beta_0 + \beta_1 X_i + \varepsilon_i$ , but we observed $X_i' = X_i + \delta_i$ . Assume $\delta_i \sim N(0, \sigma^2)$ , then we get

$Y_i + \beta_0 + \beta_1 X_i' + (\varepsilon_i - \beta_i \delta_i)$ . If $\delta$ , $\varepsilon$ , $X$ are all independent and normally distributed, it is known that $(Y, X')$ is bivariate

normal with conditional mean $E(Y|X') = \beta_0' + \beta_1' X'$ , where $\beta_0' = \beta_0$ and $\beta_1' = \dfrac{\beta_1}{1 + \sigma_\delta^2 / \sigma_X^2}$ .

- The estimated slope from regression of $Y$ on $X'$ is a biased estimate of $\beta_1$ .

- If the goal is to estimate $\beta_1$ without bias, then estimate $\dfrac{\sigma_\delta^2}{\sigma_X^2}$ from data to get an approximate unbiased estimate

    of $\beta_1$ .

- If goal is to estimate predictions for $Y$, presence of errors in $X$ decreases accuracy of predictions. A common
    assumption is to assume $\sigma_\delta^2$ small compared to $\sigma_Y^2$ and $\sigma_X^2$, then the bias of $\beta_1$ is small and reduction of
    accuracy in prediction is minimal. Then proceed as if there was no measurement error.

## SIMPLE LINEAR REGRESSION IN MATRIX TERMS

### Background

Let $\vec{X} = \begin{pmatrix} X_1 \\ \vdots \\ X_n \end{pmatrix}$ . Let $\vec{A}$ be a scalar, vector, or matrix.

Expectation: $E(\vec{X}) = \begin{pmatrix} E(X_1) \\ \vdots \\ E(X_n) \end{pmatrix}$ , $E(\vec{X} + \vec{Y}) = E(\vec{X}) + E(\vec{Y})$ , $E(\vec{A}\vec{X}) = \vec{A} E(\vec{X})$ .

Variance or Variance-Covariance Matrix: $\text{var}(\vec{X}) = \text{cov}(\vec{X}) = E[(\vec{X} - E(\vec{X}))(\vec{X} - E(\vec{X}))^T]$ is a symmetric matrix where the $i^{\text{th}}$ diagonal element is $\text{var}(X_i)$ and the $ij^{\text{th}}$ element is $\text{cov}(X_i, X_j)$ . Note that $\text{cov}(\vec{A}\vec{X}) = \vec{A}\,\text{cov}(\vec{X})\,\vec{A}^T$

### Definitions

1. $\vec{Y} = \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix}$ is a random vector ( $n \times 1$ ).

2. $\vec{\varepsilon} = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix}$ is a random vector ( $n \times 1$ ).

3. $\vec{\beta} = \begin{pmatrix} \beta_1 \\ \beta_1 \end{pmatrix}$ is not a random vector ( $2 \times 1$ ).

4. $\vec{X} = \begin{bmatrix} 1 & X_1 \\ \vdots & \vdots \\ 1 & X_n \end{bmatrix}$ is a matrix, maybe random ( $n \times 2$ ).

Then $\underset{n \times 1}{\vec{Y}} = \underset{n \times 1}{\vec{X}\vec{\beta}} + \underset{n \times 1}{\vec{\varepsilon}} = \begin{pmatrix} \beta_0 + \beta_1 X_1 + \varepsilon_1 \\ \vdots \\ \beta_0 + \beta_1 X_n + \varepsilon_n \end{pmatrix}$ .

### Gauss-Markov Assumptions

In terms of $\vec{\varepsilon}$ :

- $E(\vec{\varepsilon}) = \vec{0}$ .
- $\text{cov}(\vec{\varepsilon}) = \sigma^2 \vec{I}$ .

In terms of $\vec{Y}$ :

- $E(\vec{Y}) = X \vec{\beta}$ .
- $\text{cov}(\vec{Y}) = \sigma^2 \vec{I}$ .

Additional assumption: $\varepsilon$ has an *n*-dimensional normal distribution.

## Least Square Estimate

Find $\beta_0$ and $\beta_1$ that minimizes the sum of squares of residuals:

$$S = \sum (Y_i - \beta_0 - \beta_1 X_i)^2$$
$$= (\vec{Y} - X\beta)^T (\vec{Y} - X\vec{\beta})$$
$$\vdots$$
$$= \vec{Y}\vec{Y}^T - 2\vec{\beta}^T \vec{X}^T \vec{Y} + \vec{\beta}^T \vec{X}^T \vec{X}\vec{\beta}$$

Set $\dfrac{\partial S}{\partial \vec{\beta}} = 0$ , we get $-2\vec{X}^T \vec{Y} + 2\vec{X}^T \vec{X}\vec{\beta} = 0 \Rightarrow \vec{X}^T \vec{X}\vec{\beta} = \vec{X}^T \vec{Y} \Rightarrow \vec{\beta} = (\vec{X}^T \vec{X})^{-1} \vec{X}^T \vec{Y}$ assuming $(\vec{X}^T \vec{X})^{-1}$ exists.

Now, $\text{rank}(\vec{X}^T X) = \text{rank}(\vec{X}) = 2$ so $\vec{X}^T \vec{X}$ invertible, and

$$(\vec{X}^T \vec{X})^{-1} = \frac{1}{n \sum X_i^2 - (\sum X_i)^2} \begin{bmatrix} \sum X_i^2 & -\sum X_i \\ -\sum X_i & n \end{bmatrix} = \begin{bmatrix} \dfrac{\sum X_i^2}{n S_{XX}} & \dfrac{-\bar{X}}{S_{XX}} \\ \dfrac{-\bar{X}}{S_{XX}} & \dfrac{1}{S_{XX}} \end{bmatrix} .$$

Therefore, $\vec{b} = (\vec{X}^T X)^{-1} \vec{X}^T Y = \begin{bmatrix} \dfrac{\sum X_i^2}{n S_{XX}} & \dfrac{-\bar{X}}{S_{XX}} \\ \dfrac{-\bar{X}}{S_{XX}} & \dfrac{1}{S_{XX}} \end{bmatrix} \begin{bmatrix} \sum Y_i \\ \sum X_i Y_i \end{bmatrix} = \begin{bmatrix} \dfrac{\sum X_i^2 \sum Y_i}{n S_{XX}} - \dfrac{\bar{X} \sum X_i Y_i}{S_{XX}} \\ \dfrac{-\bar{X} \sum X_i Y_i}{S_{XX}} + \dfrac{\sum X_i Y_i}{S_{XX}} \end{bmatrix} = \begin{bmatrix} b_0 \\ b_1 \end{bmatrix}$ .

## Properties of Least Squares Estimators

Assume $X$ constant. We know $b = (X^T X)^{-1} X^T Y = AY$ . We have

- $E(b) = \beta$ ,
- $\text{cov}(b) = \sigma^2 \begin{bmatrix} \dfrac{\sum X_i^2}{S_{XX}} & \dfrac{-barX}{S_{XX}} \\ \dfrac{-barX}{S_{XX}} & \dfrac{1}{S_{XX}} \end{bmatrix}$ .

## Fitted Values and Residuals

$$e = \begin{pmatrix} e_1 \\ \vdots \\ e_n \end{pmatrix} , \quad \hat{Y} = \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix} .$$

Residuals: $e = Y - \hat{Y} = Y - Xb = (1 - H)Y$ .

Notice $\hat{Y} = Xb = X(X^T X)^T X^T Y = HY$ , where $H = X(X^T X)^T X^T$ is the "hat" matrix. $H$ is symmetric and idempotent. $\text{cov}(e) = \sigma^2 (I - H)$ . Note the df for $e^T e = \sum e_i^2 = \text{SSE}$ is $n - 2$ , which is $\text{rank}(I - H)$ .

## Analysis of Variance, Decomposition

Decomposed SSTO: $\sum (Y_i - \bar{Y})^2 = \sum Y_i^2 - n\bar{Y}^2$ . In matrix form, $\text{SSTO} = Y^T Y - \dfrac{1}{n} Y^T J Y$ where $J$ is an $n \times n$ matrix of 1's.

Now, $Y^T Y = e^T e + b^T X^T X b$ , and so

$$Y^T Y - \frac{1}{n} Y^T J Y = b^T X^T X b - \frac{1}{n} Y^T J Y + e^T e \; .$$

$$\underbrace{\phantom{Y^T Y - \frac{1}{n} Y^T J Y}}_{\text{SSTO}} \quad \underbrace{\phantom{b^T X^T X b - \frac{1}{n} Y^T J Y}}_{\text{SSR}} \quad \underbrace{\phantom{e^T e}}_{\text{SSE}}$$

Each of these is a quadratic form:

- $\text{SSTO} = Y^T \left( I - \frac{1}{n} J \right) Y$ ,
- $\text{SSR} = Y^T \left( H - \frac{1}{n} J \right) Y$ ,
- $\text{SSE} = Y^T (I - H) Y$ .

## ANOVA Table

| Source | SS | df | MS |
|--------|----|----|----|
| Regression | $b^T X^T X b - \frac{1}{n} Y^T J Y$ | 1 | $b^T X^T X b - \frac{1}{n} Y^T J Y$ |
| Error | $e^T e$ | $n-2$ | $s = \dfrac{e^T e}{n-2}$ |
| Total | $Y^T Y - \frac{1}{n} Y^T J Y$ | $n-1$ | |

## Confidence and Prediction Intervals

- The standard error for mean value of $Y$ at $X_h$ is $s \sqrt{X_h^T (X^T X)^{-1} X_h}$ where $X_h = \begin{pmatrix} 1 \\ \text{value of } X \end{pmatrix}$ .
- The standard error for predicted value of $Y$ at $X_h$ is $s \sqrt{1 + X_h^T (X^T X)^{-1} X_h}$ where $X_h = \begin{pmatrix} 1 \\ \text{value of } X \end{pmatrix}$ .

# Multiple Linear Regression

More than one predictor variable. Why?
- Often multiple $X$'s arise naturally.
- Control for some $X$'s.
- Want to fit a polynomial.
- Want to compare regression line for two groups.

## Multiple Regression Model

$k$ is the number of $X$'s.

The model: $Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \cdots + \beta_k X_{ik} + \varepsilon_i \quad i = 1, \ldots, n$ .

In matrix form: $Y = X \beta + \varepsilon$ , where $\underset{n \times 1}{Y} = \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix}$ , $\underset{(k+1) \times 1}{\beta} = \begin{pmatrix} \beta_0 \\ \vdots \\ \beta_k \end{pmatrix}$ , $\underset{n \times 1}{\varepsilon} = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix}$ , $\underset{n \times (k+1)}{X} = \begin{pmatrix} 1 & X_{11} & \cdots & X_{1k} \\ \vdots & \vdots & \vdots & \vdots \\ 1 & X_{n1} & \cdots & X_{nk} \end{pmatrix}$ .

We need to estimate $k+1$ unknown regression coefficients and 1 unknown variance, so we need $n > k + 2$ observations.

## Gauss-Markov Conditions

$$E(\varepsilon) = 0 \qquad\qquad E(Y) = X \beta$$
$$\text{cov}(\varepsilon) = \sigma^2 I \qquad\qquad \text{cov}(Y) = \sigma^2 I$$

**Least Square Estimate, Residuals, and Standard Error**

- LS estimate of $\beta$ is $b = \begin{pmatrix} b_0 \\ \vdots \\ b_k \end{pmatrix} = \left(X^T X\right)^{-1} X^T Y$ .

- Residuals is $e = Y - \hat{Y} = (I - H) Y$ .

- The estimate for $\sigma^2$ is $s^2 = \dfrac{e^T e}{\mathrm{df}} = \dfrac{e^T e}{n - k - 1}$ . This is an unbiased estimator, i.e. $E(s^2) = \sigma^2$ .

**General Comments**

- Single response, multiple explanatory variables.
- Regression gives mean response for each combination of explanatory variables.
- It won't be a helpful function if it is very complicated.
- It is usually unwise to think that there is some exact/discoverable/estimable regression equation – many possible models are available. One or two models may adequately approximate the mean of the response as a function of the explanatory variables.

**Interpreting Estimated Coefficients**

- $b_0$ : Estimate of the mean value of $Y$ when all the $X$'s are 0.
- $b_j$ : Estimated change on $Y$ due to a unit increase in $X_j$ with all other variables held constant (which may be impossible, especially in observational studies).

**Assessing Model Assumptions**

Use residual plots as in simple linear regression.

- Residuals vs. predicted value: non-constant variance, outliers, influential points.
- Residuals vs. $X_j$ : curvature.

# PARAMETER ESTIMATE

| Variable | df | Parameter Estimate | Standard Error | t-value | Pr > \|t\| |
|---|---|---|---|---|---|
| Intercept | 1 | $b_0$ | $\mathrm{se}(b_0)$ | $\dfrac{b_0}{\mathrm{se}(b_0)}$ | p-value from $t_{n-k+1}$ distribution |
| $X_1$ | 1 | $b_1$ | $\mathrm{se}(b_1)$ | $\dfrac{b_0}{\mathrm{se}(b_1)}$ | p-value from $t_{n-k+1}$ distribution |
| ...etc... | ...etc... | ...etc... | ...etc... | ...etc... | |
| $X_k$ | 1 | $b_k$ | $\mathrm{se}(b_k)$ | $\dfrac{b_0}{\mathrm{se}(b_k)}$ | p-value from $t_{n-k+1}$ distribution |

Notes:

- The estimate of $\mathrm{var}(b_j)$ is the $j$-th diagonal entry of $s^2 \left(X^T X\right)^{-1} = \mathrm{MSE} \left(X^T X\right)^{-1}$ , and $\mathrm{se}(b_j) = \sqrt{\mathrm{var}(b_j)}$ .

- To test $H_0 : \beta_j = 0$ vs $H_A : \beta_j \neq 0$ , use test statistic $t_{\mathrm{obs}} = \dfrac{b_j}{\mathrm{se}(b_j)}$ , which will be from a value from $t_{n-k+1}$ distribution under $H_0$ . This test gives an indication of whether or not the $j$-th predictor variable statistically significantly contributes to the prediction of the response variable over and above all of the other predictor variables.

- The CI for $\beta_j$ is $b_j \pm t_{n-k-1, 1-\alpha/2} \, \mathrm{se}(b_j)$ . Use Bonferroni if constructing CI for multiple $\beta_j$'s.

## ANOVA TABLE

| Source | df | SS | MS | F | Pr > F |
|---|---|---|---|---|---|
| Regression | $k$ | $Y^T\left(H-\dfrac{1}{n}J\right)Y$ | $MSR=\dfrac{SSR}{k}$ | $\dfrac{MSR}{MSE}$ | |
| Error | $n-k-1$ | $Y^T(I-H)Y$ | $MSE=\dfrac{SSE}{n-k-1}$ | | |
| Total | $n-1$ | | | | |

### F-Test From ANOVA

$H_0: \beta_0=\beta_1=\cdots=\beta_k=0$     vs.   $H_A:$ at least one $\beta_j\neq 0$ .

Test statistic: $F_{obs}=\dfrac{MSR}{MSE}$ . Under $H_0$, $F_{obs}$ is an observation from a $F_{k,n-k-1}$ distribution.

This test answers "Does regression do anything?".

# R² AND ADJUSTED R²

### Coefficient of Multiple Determination

$R^2=\dfrac{SSR}{SSTO}$ is the percentage of variation in the $Y$'s explained by the regression equation. It is called the coefficient of multiple determination.

It is not the square of a correlation coefficient and $\sqrt{R^2}$ is not really useful in multiple regression.

In multiple regression, $R^2$ is not that great of a judge of the model because the more predictor variables the greater $R^2$, despite whether or not the predictor variables are useful for estimating the response variable.

### Adjusted Coefficient of Multiple Determination

An attempt to make $R^2$ more useful is to use adjusted $R^2$. It is adjusted for the number of predictor variables in the model.

$\text{adj } R^2=1-(n-1)\dfrac{MSE}{SSTO}=1-\dfrac{n-1}{n-k-1}\dfrac{SSE}{SSTO}$ .

It increases only if MSE increases. It is useful for choosing the best model.

## CONFIDENCE AND PREDICTION INTERVALS FOR MEAN RESPONSE

Let $X_h=\begin{pmatrix}1\\X_{h1}\\\vdots\\X_{hk}\end{pmatrix}$ .

The standard error for $E(Y)$ at $X_h$ is $s\sqrt{X_h^T(X^TX)^{-1}X_h}$ .

The standard error for the predicted value of $Y$ is $s\sqrt{1+X_h^T(X^TX)^{-1}X_h}$ .

## INDICATOR/DUMMY VARIABLES

For a categorical variable with $j$ levels, define $j-1$ variables $I_i=\begin{cases}1 & \text{if in the i-th level}\\0 & \text{otherwise}\end{cases}$ , $i=1,2,\ldots,j-1$ .

**Advantages vs Disadvantages**

| Advantages of Indicator/Dummy Variables | Disadvantages of Indicator/Dummy Variables |
|---|---|
| • Useful if we expect a complex relationship between predictors and the response variable.<br>• "Non-parametric" approach – doesn't assume form for functional relationship. | 1. Lose functional form of relationship.<br>2. Lose ordering.<br>3. Costs df. |

# PARTIAL F-TEST

To test whether a subset of the $\beta$ 's is 0.

We have

- SSE in reduced model > SSE in full model,
- SSR in reduced model < SSR in full model,
- SSTO in reduced model = SSTO in full model,

It can be shown that $F_{\text{obs}} = \dfrac{(\text{SSE}_{\text{reduced}} - \text{SSE}_{\text{full}})/(\text{df}_{\text{error, reduced}} - \text{df}_{\text{error, full}})}{(\text{SSE}_{\text{full}})/(\text{df}_{\text{error, full}})}$ is an observation from a $F_{\text{df}_{\text{error, reduced}} - \text{df}_{\text{error, full}}, \text{df}_{\text{error, full}}}$

distribution.

# PARTIAL CORRELATION

Use when considering the reduced vs. full model where the full model has one additional predictor variable. It is a measure of the contribution of the additional predictor variable given the others in the model.

The coefficient of partial correlation is $\sqrt{\dfrac{\text{SSE}_{\text{reduced}} - \text{SSE}_{\text{full}}}{\text{SSE}_{\text{reduced}}}}$ . It is negative if the coefficient of the additional predictor variable is negative, otherwise it is positive.

# POLYNOMIAL REGRESSION

## Second Order (Quadratic)

A second order (quadratic) model in one variable: $Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \varepsilon$ .

Full second order model with two variables: $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_{11} X_1^2 + \beta_{22} X_2^2 + \beta_{12} X_1 X_2 + \varepsilon$ .

Notice that it is still linear regression.

## Comments

- As a general rule, when higher order terms are in the model, also include lower order terms in that variable, even if their coefficients not statistically different from 0.
- Idea: Higher order terms are refinements to base effect of predictor variables on $Y$ .
- The meaning of individual coefficients are difficult and not necessary to interpret; can't change $X$ and hold $X^2$ constant.

## Multi Collinearity and Centering

A problem that often arises when doing multiple regression is that $X$ and $X^2$ are often highly correlated. When predictor variables are highly correlated, the fitted equation is unstable. The estimated coefficients can change dramatically when data changes a little and when the set of predictor variables that are in the model changes. This is multi collinearity.

In the polynomial setting, often improves the situation to center the $X$ 's, i.e. use $\tilde{X}_{i1} = X_{i1} - \bar{X}_1$ (where $\bar{X}_1 = \dfrac{1}{n} \sum_{i=1}^{n} X_{i1}$ )

in the model instead of $X_{i1}$ .

In matrix form, $Y = \tilde{X}\tilde{\beta} + \varepsilon$ , where $\tilde{\beta} = \begin{pmatrix} \gamma_0 \\ \beta_1 \end{pmatrix}$ , $\tilde{X} = \begin{bmatrix} 1 & X_{11} - \bar{X}_1 \\ \vdots & \vdots \\ 1 & X_{n1} - \bar{X}_1 \end{bmatrix}$ . Then $\tilde{X}^T \tilde{X} = \begin{bmatrix} n & 0 \\ 0 & \sum (X_{i1} - \bar{X}_1)^2 \end{bmatrix}$ is easy to invert.

For multiple regression, $\underbrace{\tilde{X}^T \tilde{X}}_{(k+1) \times (k+1)} = \begin{bmatrix} n & 0 & \cdots & 0 \\ 0 & \ddots & \vdots & \ddots \\ \vdots & & Z^T Z & \\ 0 & \ddots & \vdots & \ddots \end{bmatrix}$ where $\left[ Z^T Z_{jj} \right] = \sum_{i=1}^{n} (X_{ij} - \bar{X}_j)^2$ (diagonal) and

$\left[ Z^T Z_{jk} \right] = \sum_{i=1}^{n} (X_{ij} - \bar{X}_j)(X_{ik} - \bar{X}_k)$ $j \neq k$ (close to 0 if little correlation).

## INTERACTION

### Definition
Two explanatory variables are said to interact if the effect that one has on the mean response depends on the value of the other.

### Interaction Term
To model interaction, introduce a new variable $X_i X_j$ in model.

### When To Add Interaction Terms
- Question of interest has to do with interaction.
- Suspect interaction exists because of the nature of what's being measured or because plot of residual vs. interaction shows a relationship.
- If an interaction term for 2 predictor variables is in model, should include terms for predictor variables individual (i.e. if $X_1 X_2$ is in model, then include $X_1$ and $X_2$ ).

## MULTICOLLINEARITY

### Effects of Multicollinearity
Serious correlation among predictors may have the following effects:
- Regression coefficients will change dramatically according to which other variables are included/excluded from the model. In worst cases, regression coefficients will be dramatically large with signs that may be opposite to what you expect.
- Standard error of regression coefficient will be large.
- Predictors that are known to be related with $Y$ will not have their coefficients be statistically different from 0.

### Checking for Multicollinearity
- Check the correlations among pairs of $X$'s.
- Check if regression coefficients have wrong signs given prior knowledge.
- Check if removing or adding a predictor variable produces surprising changes in fitted model.
- Check if predictors anticipated to be important have non-significant coefficients.
- Examine $VIF$ (variance inflation factor): $VIF_i = (1 - R_i^2)^{-1}$ where $R_i^2$ is the coefficient of multiple determination obtained when the $i$-th predictor variable is regressed against other predictor variables. Large $VIF$ is a sign of multicollinearity; $VIF > 10$ deserves attention (suggestion).

- Examine tolerance $\frac{1}{VIF}$ .

## VARIABLE SELECTION

Why eliminate possible explanatory varaibles?
1. Simpler models are better than complex ones – "parsimony". Want a simple model that doesn't leave out anything important.
2. Unnecessary terms in the model result in less precise inferences ( $MSE \uparrow$ ).

### Automated Procedures

Forward Selection, Backward Elimination, Stepwise – each give one final model (but not necessarily the same).
More modern: All Possible Subsets Regression – $k$ predictor variables, $2^k$ possible models.

### Criteria for "Best"

- Adjusted $R^2$ .
- Mallow's $C_p$ : In a model with $p-1 \leq k-1$ predictor variables (so $p$ and $k$ $\beta$ 's), $C_p = \frac{SSE_p}{MSE_k} - (n - 2p)$ . For a good model, want $C_p = p$ .
- Prediction SS: Suppose the $i$ -th data point is deleted. Estimate the line without it and predict $Y_i$ for all $n$ points; call these $\hat{Y}_{i(i)}$ . Then $PRESS_p = \sum_{i=1}^{n} \left( Y_i - \hat{Y}_{i(i)} \right)^2$ . Smaller $PRESS_p$ is better. (Note: $p = k + 1$ )

## IDENTIFYING OUTLIERS AND INFLUENTIAL POINTS

### Outliers

Outliers have large residuals. How large is too large?

- Standardized residual: $\frac{e_i}{s\sqrt{1 - h_{ii}}}$ where $h_{ii}$ is the $i$ -th diagonal entry of $H$ . If $> 3$, unusual; if $> 4$, very unusual.
- Studentized residual: $e_i^{\iota} = \frac{e_i}{s_{(i)}\sqrt{1 - h_{ii}}}$ where $s_{(i)}^2$ is the estimate of $\sigma^2$ from regression without the $i$ -th data point.
- Alternate approach: Add to regression a dummy variable $d_i = \begin{cases} 1 & \text{for } i\text{-th observation} \\ 0 & \text{otherwise} \end{cases}$ . If the coefficient of $d_i$ is statistically different from 0, then the $i$ -th data point does not fit model.

### Influential Point

An influential point pulls the lines towards itself.

- Look at $b - b_{(i)}$ where $b_{(i)}$ is the LS estimate with the $i$ -th observation omitted.
- Loot at $\hat{Y}_i - \hat{Y}_{(i)}$ .
- Cooks Distance: Looks at all predicted values when the $i$ -th observation is deleted.
- Robust Regression: For example, use Least Median Squares (LMS) which minimizes the median of square residuals. This is computationally extensive.
- Reweighed LS: Remove outliers from LMS and then do Ordinary Least Squares (OLS).