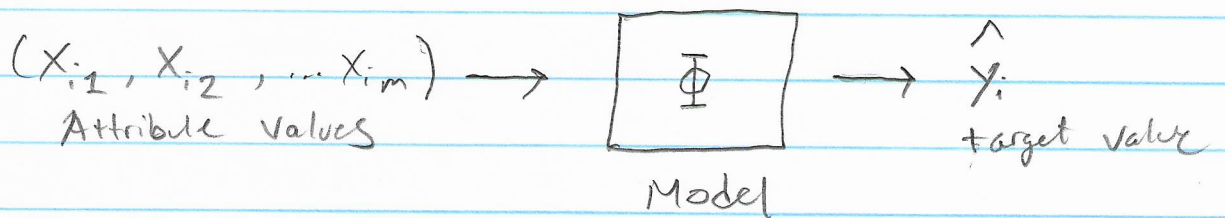TIM 245 Lecture 7 (4/24/17)

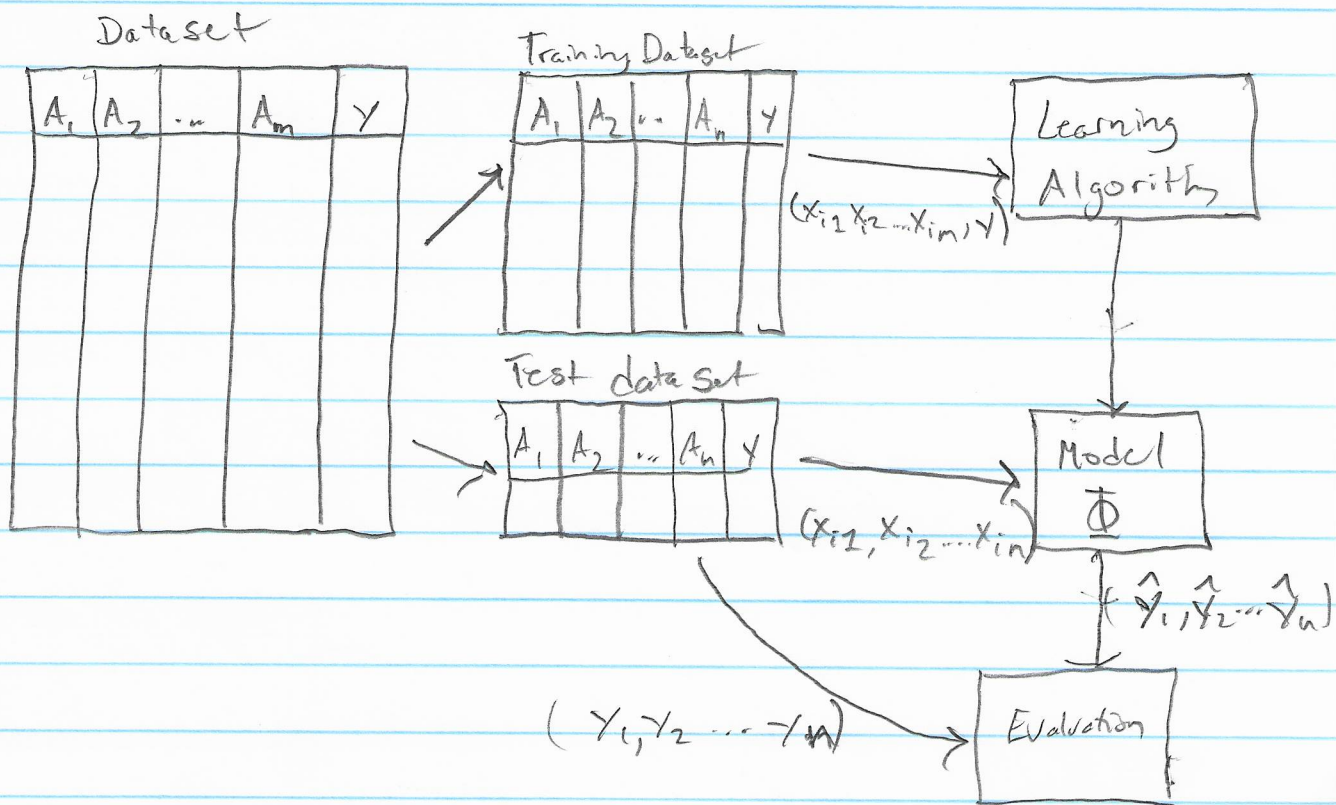<u>Agenda</u>

1) General process for supervised learning problems

2) Evaluation of classification and prediction models

3) Linear regression and introduction to prediction models

4) Work on project (Time permitting)
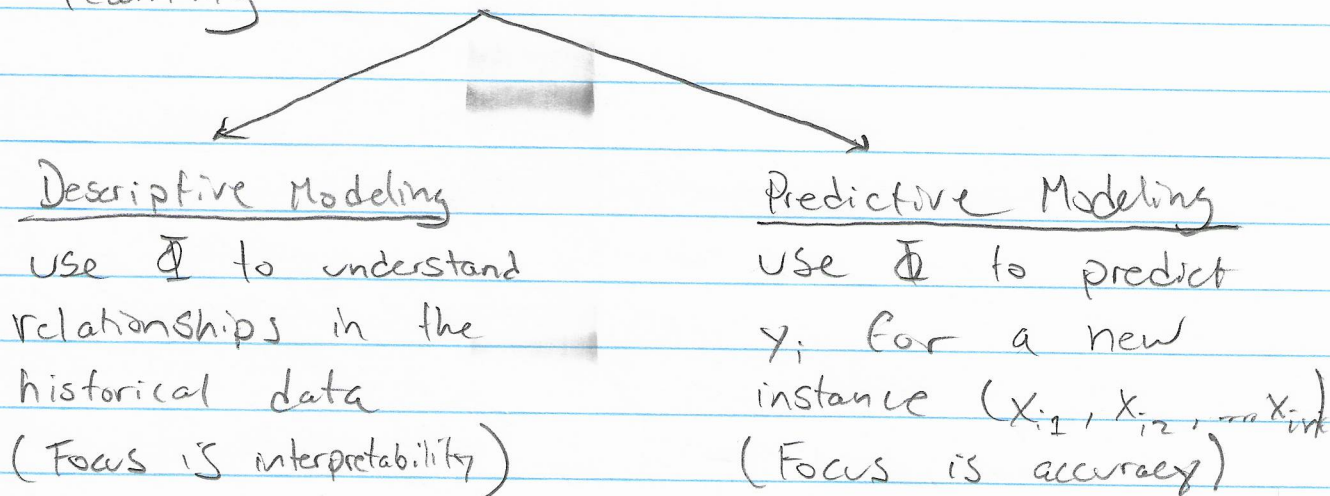
① General Process for Supervised learning problems

Objective: Create a function or model that can map an input set of attributes to a target output value

$$(X_{i1}, X_{i2}, \ldots X_{im}) \longrightarrow \boxed{\Phi} \longrightarrow \hat{Y}_i$$

Attribute Values $\qquad$ Model $\qquad$ target value

The model, $\overline{\Phi}$, is created through a supervised learning process:

Dataset

| $A_1$ | $A_2$ | ... | $A_m$ | $Y$ |
|-------|-------|-----|-------|-----|
|       |       |     |       |     |

Training Dataset

| $A_1$ | $A_2$ | ... | $A_m$ | $Y$ |
|-------|-------|-----|-------|-----|
|       |       |     |       |     |

$(X_{i1} X_{i2} \ldots X_{im}, Y)$ → Learning Algorithm

Test data set

| $A_1$ | $A_2$ | ... | $A_n$ | $Y$ |
|-------|-------|-----|-------|-----|
|       |       |     |       |     |

$(X_{i1}, X_{i2} \ldots X_{in})$ → Model $\Phi$

$\{\hat{Y}_1, \hat{Y}_2 \ldots \hat{Y}_W\}$

$(Y_1, Y_2 \ldots Y_W)$ → Evaluation

Two general applications of supervised learning

Descriptive Modeling
use $\Phi$ to understand relationships in the historical data
( Focus is interpretability )

Predictive Modeling
use $\Phi$ to predict $y_i$ for a new instance $(x_{i_1}, x_{i_2}, \ldots x_{ik})$
( Focus is accuracy )

Key issue in Descriptive Modeling : Correlation vs Causation

Four possible cases for $\Phi : X \to Y$

1) Causal Link

$$\left. \begin{array}{c} X \to Y \\ Y \to X \end{array} \right\} \text{ No sure which one}$$

2) Hidden Cause

$$\left. \begin{array}{c} Z \\ \swarrow \searrow \\ X \qquad Y \end{array} \right\} \begin{array}{l} Z \text{ is not in the} \\ \text{data set} \end{array}$$

3) Confounding Factor

$$Z \searrow$$
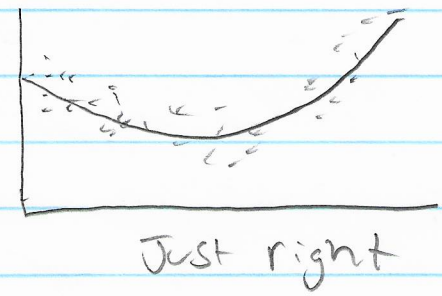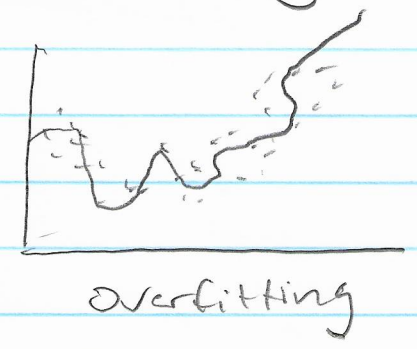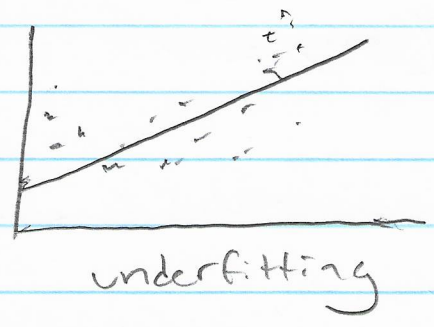$$\quad \quad Y$$
$$X \nearrow$$

4) Coincidence (Noise)

$$X \quad Y$$

Typically addressed through:
- Careful with attribute Selection
- Consulting with Subject matter experts
- Randomized experiments (A/B testing)

Key issue in Predictive Modeling : overfitting vs underfitting



underfitting



overfitting



Just right

Typically addressed through:
- Model regularization (Lasso, Ridge)
- Careful evaluation of the accuracy
- Cross validation
- Ensemble methods (Bagging, Boosting)

② Evaluation of Classification and prediction models

Classification model are evaluated using a Confusion matrix

|  |  | Predicted | |
|---|---|---|---|
|  |  | Award | No Award |
| Actual | Award | True Positives | False Negative |
|  | No Award | False Positive | True Negatives |

Measures: Accuracy, Precision, Recall, f-measure, g-mean

Example: $\text{Accuracy} = \dfrac{TP + TN}{TP + TN + FP + FN}$

Prediction Models are evaluated based on the error

$$e_i = y_i - \hat{y}_i \quad (\text{residual})$$

Measures: Sum of Squared Error (SSE), Mean Absolute Error (MAE), Mean Average Percent Error (MAPE)

Example: $SSE = \sum\limits_{i=n}^{n} (y_i - \hat{y}_i)^2$

③ Linear Regression

three general types of relationships between X and Y

<u>Positive Correlation</u>: values move in same direction

$\uparrow X \uparrow Y$    or    $\downarrow X \downarrow Y$

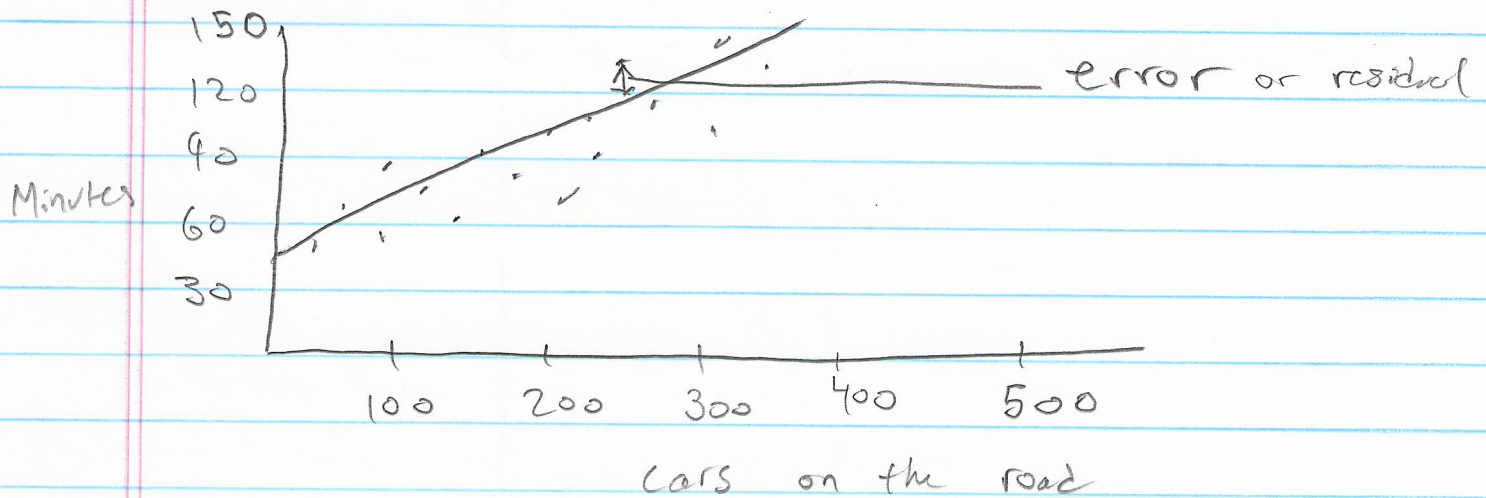<u>Negative Correlation</u>: Values move in opposite directions

$\uparrow X \downarrow Y$    or    $\downarrow X \uparrow Y$

<u>No Correlation</u>

$\uparrow X \, Y$   or   $\downarrow X \, Y$   or   $X \uparrow Y$   or   $X \downarrow Y$

Linear regression estimates the type and strength of linear relationships in the training data set, i.e. best straight line fit

Example :  HWY 17 Commute times vs number of cars on the road



General form of a linear regression model

$$\hat{y} = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_m X_{im}$$

dependent variable    intercept    Coefficients    independent variables

Coefficients indicate the strength and type of the relationship:

     Magnitude : Strength of the relationship

          $\approx 0$ : no correlation

     Sign       : type of relationship
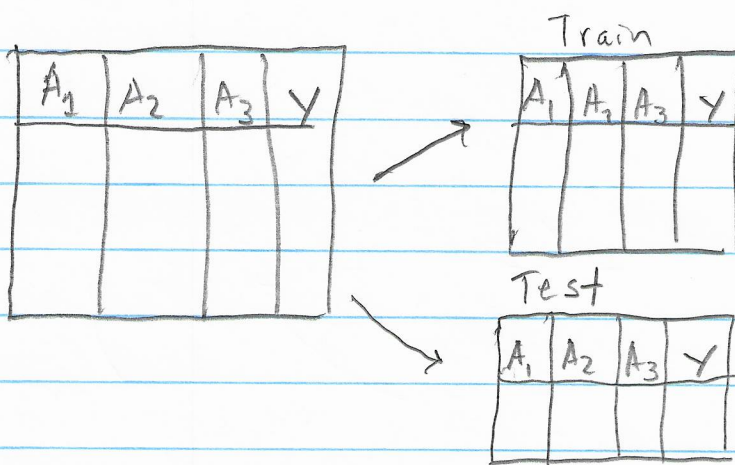
          + : positively correlated

          − : negatively correlated

Simple Linear Regression: linear regression with only a single independent variable or attribute
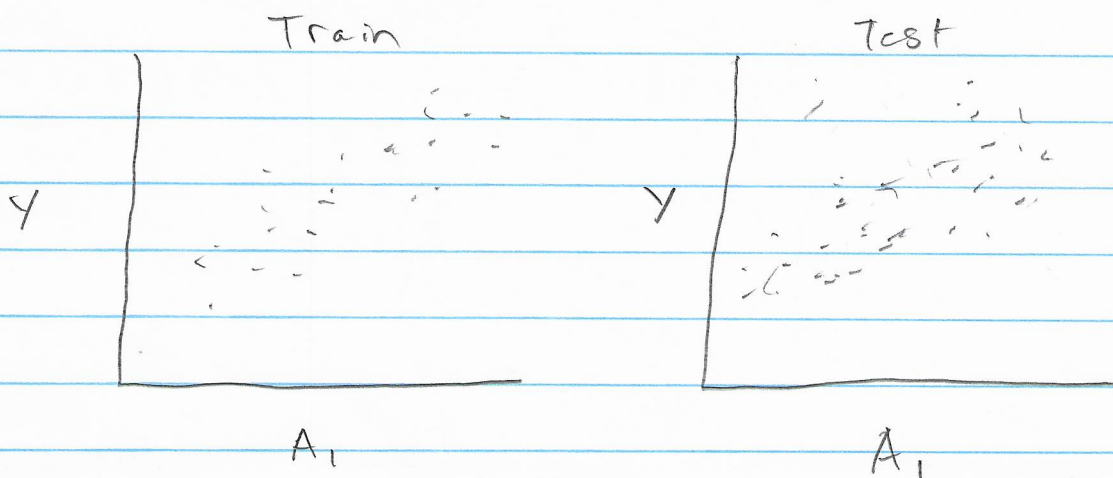
$$\hat{y} = B_0 + B_j X_{ij} \qquad (\text{think } y = mx + b)$$

Process for Simple linear Regression:

1) Separate the dataset into training dataset and test dataset



2) Plot data for $j=1$ ($A_1$)

3) Estimate $\beta_0$ (intercept) and $\beta_1$ (slope) on the training dataset for $j=1$ ($A_1$) using least squares

$$\beta_1 = \frac{\sum_{i=1}^{n} (x_{i1} - \bar{A}_1)(y_i - \bar{y})}{\sum_{i=1}^{n} (x_{i1} - \bar{A}_1)^2}$$

$$\beta_0 = \bar{y} - \beta_1 \bar{A}_1$$

where $\bar{A}_1$ and $\bar{y}$ are the means of $A_1$ and $y$, respectively

4) Calculate the sum of squared error on the test dataset

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

5) Repeat for each attribute $A_j$ in the data set ($j = 2, 3, \dots m$)

6) Select the model with the lowest error

Provides a good baseline for more sophisticated linear models