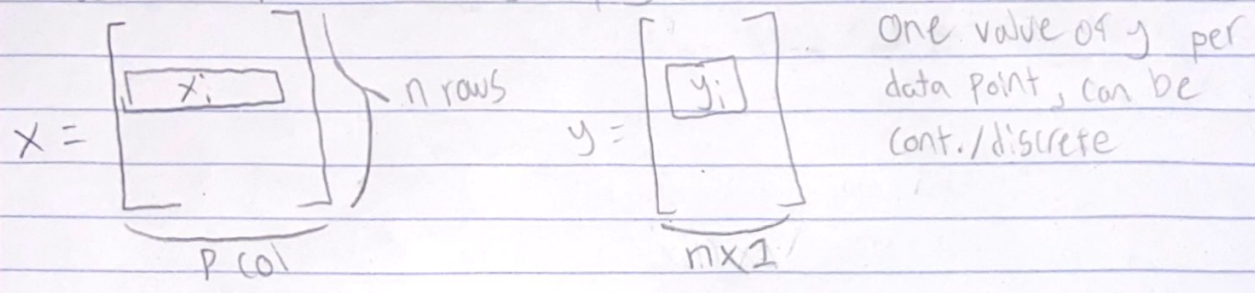# 9/9 Lecture Notes.

▷ <u>Admin</u>: Change up seats, Lab 1 due tmr

▷

▷ <u>Data repr. and featurization</u>

▷ Features (input) ex temp

▷ Label (output) ex will you play outside.

$$X = \begin{bmatrix} \boxed{x_i} \\ \ \\ \ \end{bmatrix} \Bigg\} \; n \text{ rows} \qquad y = \begin{bmatrix} \boxed{y_i} \\ \ \end{bmatrix}$$

One value of $y$ per data point, can be cont./discrete

$\underbrace{\phantom{xxxxx}}_{p \text{ col}}$ $\qquad\qquad\qquad \underbrace{\phantom{xxxx}}_{n \times 1}$

▷ Feature: name → shape
  values → circle
  vector → $X = [x_1, x_2, \cdots x_p]$

▷ <u>Featurization</u> (make numerical)  ex  False → 0
  ↳ so comp. can interpret them    True → 1

▷ More ex. on slides    Can map spectrum to values ex sunny → 1
▷ This process is done by Data Scientist    rainy → 0
▷ to suit the model best.

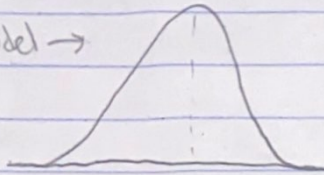▷ <u>Model</u>: distribution that captures data
▷ ex decision tree (on slide)

Params $\Bigg[$
  weather — feature
  sun ╱  ╲ rain — feat. value
  $\boxed{Y}$   $\boxed{N}$ — labels
  [1,2]   [2,0] — label counts
  2/3   2/2  — Accuracy
         (80% overall)

| Data | weather | tennis |
|------|---------|--------|
|  | S | Y |
|  | r | N |
|  | r | N |
|  | S | Y |
|  | S | N |

▷ Normal distribution model

▷ model →



$\left.\begin{array}{l}\text{mean} \\ \text{Variance}\end{array}\right]$ Model Params.

▷ Linear models:



→ Model (slope = m)    $y = mx + b$

prediction                  $m, b$ → model params

based on model

▷ Handout thoughts: Q2 - can featurize by binarizing (turning to 1s and 0s)
          Q4 - Model is "perfectly classified" depending on meaning
              ↳ if all given data points are classified
              ↳ if all possible data values are classified

▷ Purpose of models:
▷ Make predictions based on data
▷ Vizualize data
▷ Help humans make choices based on data
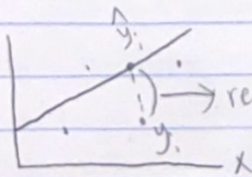▷ ↳ ex, can see that a certain feature alligns more with output

▷ Linear models:
▷ feature $\vec{x}$     eq: $h_{\vec{w}}(x) = w_0 + w_1 x = \hat{y}$     residual: $y_i - \hat{y}_i$
▷ output $y$                                prediction
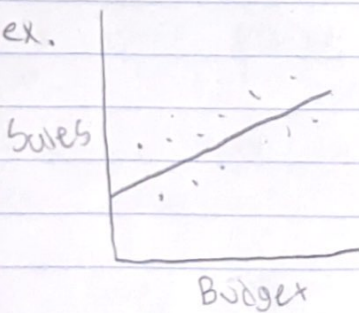          want to minimize RSS $\sum\limits_{i=1}^{?} (y_i - \hat{y}_i)^2$
          (error variance)



→ residual $(y_i - \hat{y}_i)$ → difference between prediction and data

\*\*\*

▷ Next lecture:
▷ Use CL Args to input data
▷ ↳ Parse_Args package
▷ Demo on website
▷ Notes: can use — in arg name to make it optional
▷ args are read in as string
▷ Can set default arg value
▷
▷ Goals of fitting a linear model: 1) which features impact y
2) association between x and y

▷ ex.



Sales / Budget scatter plot with fitted line → Model tells us as sales goes up, so does budget

▷ Linear Regression
▷ Used when output is continuous
▷ Learned model: linear function mapping x to y (including feature weights and bias)
▷ Goal: to minimize RSS/SSE (error)

model params

▷ "simple" lin. regress. → one feature
▷ $h_{\vec{w}} = w_0 + w_1 x = \hat{y}$

weights

▷ $\vec{w} = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}$

▷ Linear model may not always be most accurate
▷ ↳ could use polynomial curve to better fit data
▷ MSE (average error)

▷ Performance may vary from training data to testing data
▷ "Overfitting" - over performs on training data, underperforms on new data