

Overview of ML

- Train,
- Test,
- Performance

$N=1000$, $y = 1$ or 0 (target)
 x_1 - Age
 x_2 - \$

$y = \hat{f}(x_1, x_2)$
 y - true party
 \hat{y} - predicted party
 $\hat{f}(\cdot)$

$$E(y_i - \hat{y}_i)^2 \Rightarrow \arg \min$$

$$\arg \min_{\hat{f}(x_1, x_2)} E(y_i - \hat{y}_i)^2$$

$$\Rightarrow \text{ " } E(y_i - \underline{\underline{\hat{f}(x_1, x_2)}})^2$$

① Training

② Performance

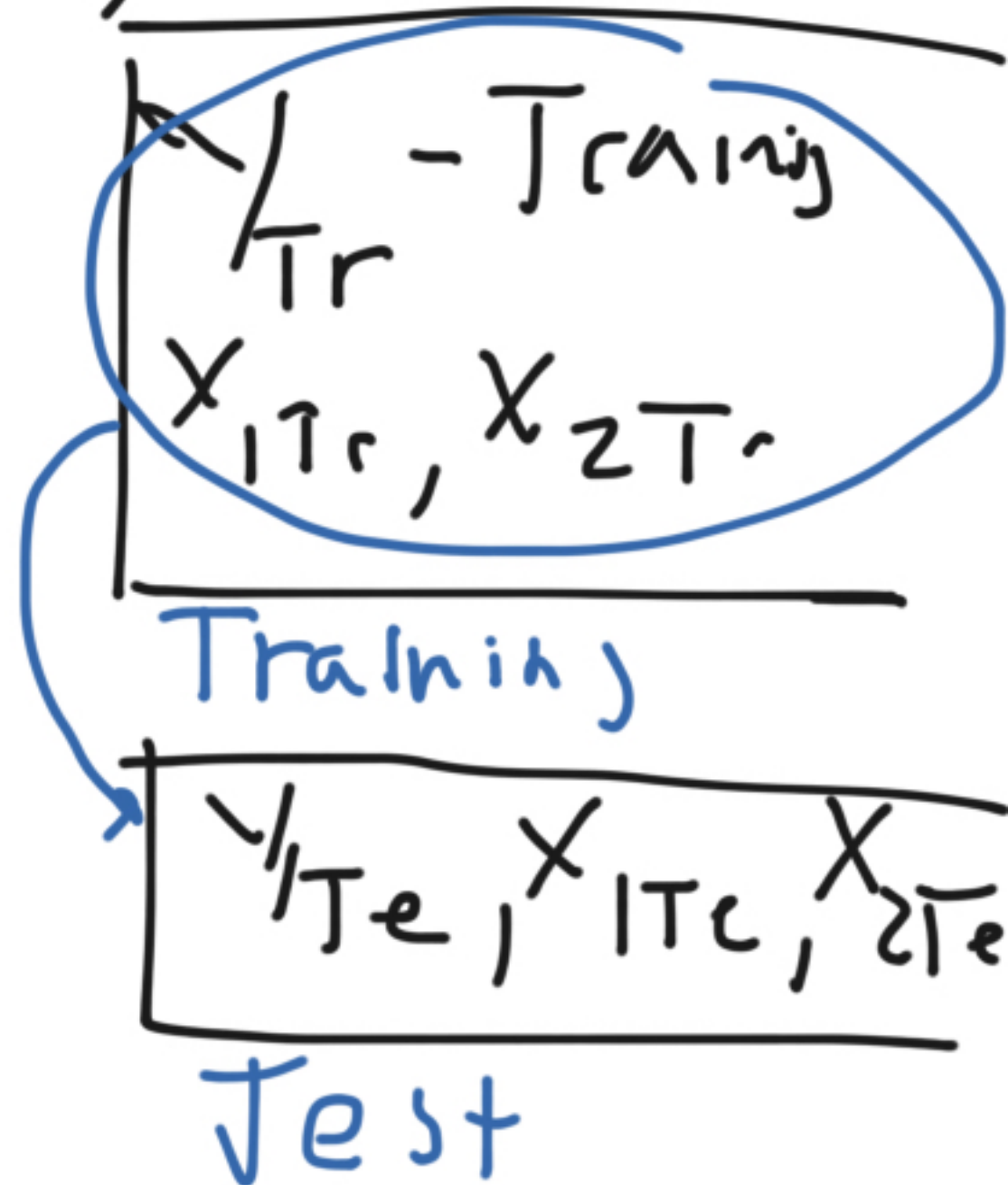
SI Training

$N=1000$ y , Age, Income

Training / Test
80% / 20%
(800) / (200)

$$y_{Tr} = \hat{f}(x_{1Tr}, x_{2Tr}) \quad ①$$

$$\hat{y}_{Te} = \hat{f}(x_{1Te}, x_{2Te})$$



$$y_{Te} - \hat{y}_{Te}$$

Divide the data? 80/20
to min bias/max representation
of the training data

<u>Person</u>	<u>Party</u>	<u>Age</u>	<u>\$</u>	<u>Person R</u>
Tr {	1	} Trenton		3
	0			101
	0			20
	1			.
Te {	.	} Princeton	800	.
	.			.
	.			.
	1000			.

Randomly select 800 obs for
training

✓ (1) $y_{Tr} = \hat{f}(x_{1Tr}, x_{2Tr})$ N_{Tr}
800

✓ (2) $\hat{y}_{Te} = \hat{f}(x_{1Te}, x_{2Te})$ N_{Te}
200

✓ (3) Compare

\hat{y}_{Te}	y_{Te}
Pred	True

Step 2 - Performance Metrics

→ classification - categorical labels

→ regression - continuous value
(prices, % dem vol)

discrete {
2 - I
1 - b
0 - R

discrete

% dem (0 - 100%)

Continuous

discrete {
2 - } # of
1 - } people in a
0 - } room

$$Acc = \boxed{80\%}$$

$$Precision = \boxed{90\%}$$

$$Recall (sensitivity) = \frac{TP}{TP + FN} = \frac{90}{120} = \boxed{75\%}$$

$$F_1 = 2 \times \left(\frac{Precision \times Recall}{Precision + Recall} \right) = 2 \left(\frac{0.90 \cdot 0.75}{0.9 + 0.75} \right)$$

$$= \boxed{0.81}$$

$$\underline{Specificity} = \frac{TN}{TN + FP}$$

True negative
rate

Performance (regression)

$$RMSE = \sqrt{\frac{\sum_{i=1}^N (y_i - \hat{y}_i)^2}{N}}$$

$N=1000$ cities

$y = \% \text{ voting for HRC in each city}$

$x_1 = \text{Age (avg for each city)}$

$x_2 = \$ \text{ (avg for each city)}$

80/20 Training / Test

Train

$$y_{Tr} = \alpha + \beta_1 x_{1Tr} + \beta_2 x_{2Tr}$$

Test

$$\hat{y}_{Te} = \alpha + \beta_1 x_{1Te} + \beta_2 x_{2Te}$$

y_{Te} - true % voting for
HRC in test data

\hat{y}_{Te} - predicted y_0 in
" " " "

$$\text{RMSE}_{T_e}$$

$$N_e = 200$$

$$\sqrt{\frac{\sum_{i=1}^{N_e} (y_i - \hat{y}_i)^2}{N_e}}$$

Bias-Variance Tradeoff

$$\underline{E(y - \hat{y})} = \text{Var}(\hat{f}(x)) + \underbrace{\text{Bias}(\hat{f}(x))}_{\downarrow}$$

$$\hat{y} = \hat{f}(x)$$

How well
the algorithm
predicts on
other datasets

(RMSE_{Te})

Predictions
on training
data

(RMSE_{Tr})

RMSE

