

SIS - III

Conceptual questions:

1.

- Binary classification: water on/off, spam/not spam
- Multiclass: blood type(A/B/O/...); digits (0,1,2,3,...)
- Regression: predict weather / price for something

y_i	1	1	0	0	
\hat{y}_i	0.9	0.9	0.3	0.6	

(using base-10 logs)

Data: $(y, \hat{y}) = (1, 0.9), (1, 0.4), (0, 0.3), (0, 0.6)$

Formula: $J(\theta) = -\frac{1}{n} \sum_{i=1}^n [y_i \ln(\hat{y}_i) + (1-y_i) \ln(1-\hat{y}_i)]$

$$= \frac{1}{4} \left[(-1 \cdot \ln(0.9) + 0 \cdot \ln(0.1)) + (-(\ln(0.4) + 0 \cdot \ln(0.6))) + \right. \\ \left. + (- (0 \cdot \ln(0.3) + 1 \cdot \ln(0.7)) + (-(\theta \cdot \ln(0.6) + 1 \cdot \ln(0.4))) \right] = \\ = \frac{1}{4} \left(-\ln(0.9) - \ln(0.4) - \ln(0.7) - \ln(0.4) \right) = \frac{2.2946}{4} \approx 0.5737$$

nuts

3. we notice that the model is overfitting

4. b

5. b

6. c

7.

Error rate = $1 - \frac{n_i}{n}$, it's minimal when ℓ is the majority class: $\ell = \arg \max_k n_k$.

Minimum error: $\gamma = \frac{\max_k n_k}{n}$

8. $f(x_c) = c$

$$J(c) = \frac{1}{n} \sum_{i=1}^n (y_i - c)^2$$

$$\frac{\partial J}{\partial c} = \frac{1}{n} \sum_{i=1}^n 2(c - y_i) = \frac{2}{n} \left(nc - \sum_{i=1}^n y_i \right)$$

1st.

$$\frac{2}{n} \left(nc - \sum_{i=1}^n y_i \right) = 0$$

$$\text{2nd} \quad \frac{d^2 J}{dc^2} = \frac{2}{n} \sum_{i=1}^n 1 = \frac{2}{n} \cdot n = 2 > 0$$

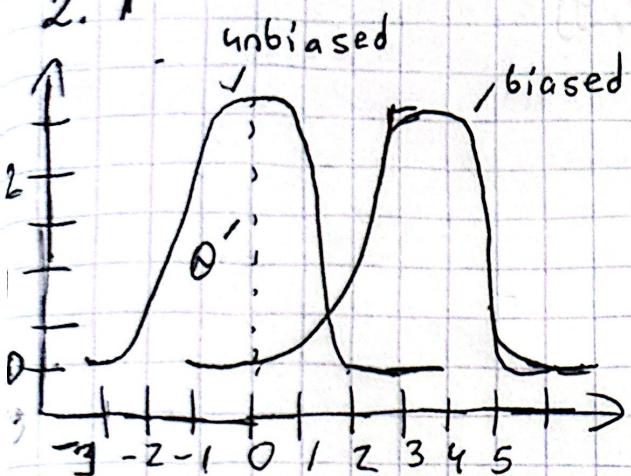
$$nc = \sum_{i=1}^n y_i$$

$$\text{so, } c = \bar{y} = \frac{1}{n} \sum_{i=0}^n y_i$$

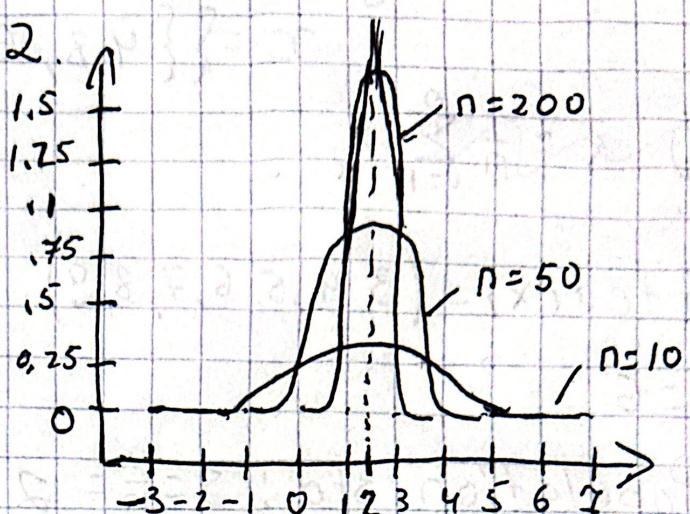
$$c = \frac{1}{n} \sum_{i=1}^n y_i = \bar{y}$$

Properties of an Estimator

2.1



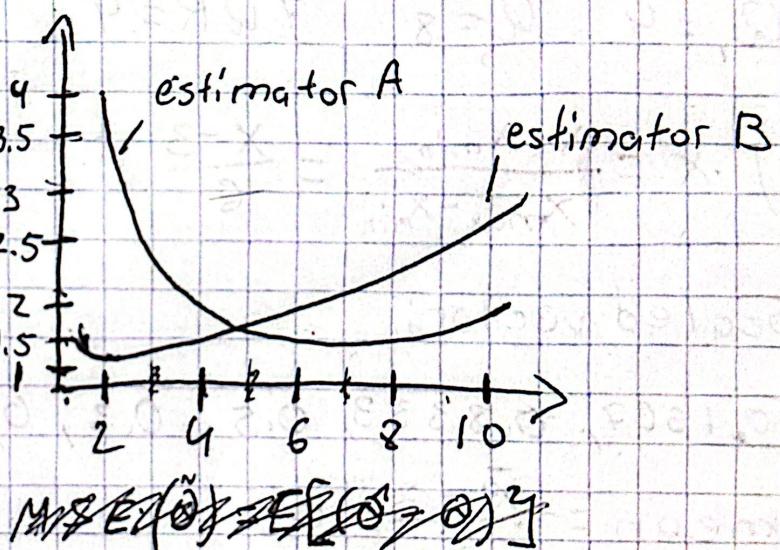
2.2



2.3



2.4



$$MSE(\hat{\theta}) = E[(\hat{\theta} - \theta)^2]$$

$$= E[(\hat{\theta} - E(\hat{\theta}) + E(\hat{\theta}) - \theta)^2]$$

$$= E[(\hat{\theta} - E(\hat{\theta}))^2] + (E(\hat{\theta}) - \theta)^2 + 2(E(\hat{\theta}) - \theta) E(\hat{\theta} - E(\hat{\theta}))$$

$$\text{Var}(\hat{\theta})$$

$$\text{Bias}(\hat{\theta}, \theta)^2$$

~~$$= E[(\hat{\theta} - E(\hat{\theta}))^2] + \text{Var}(\hat{\theta}) + \text{Bias}(\hat{\theta}, \theta)^2$$~~

$$= \text{Var}(\hat{\theta}) + \text{Bias}(\hat{\theta}, \theta)^2 - \text{systematic error}$$

how much estimator fluctuates across diff. samples

Data Scaling:

$$x = \{4, 8, 6, 3, 3, 7, 9\}$$

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\text{sorted}(x) = (3, 4, 5, 6, 7, 8, 9)$$

$$\bar{x} = 6$$

$$\cdot \text{Population std: } s = \sqrt{\frac{1}{n}} = 2$$

$$\text{Median} = 6$$

$$Q_1 = 4 \quad Q_3 = 8 \quad IQR = 4$$

$$a) \cdot x' = \frac{x - x_{\min}}{x_{\max} - x_{\min}} = \frac{x - 3}{6}$$

Scaled vector:

$$[0.1667, 0.8333, 0.5, 0.3, 0, 0.667, 1]$$

$$\text{mean} = \frac{\bar{x} - 3}{6} = 0.5$$

$$\text{std} = \frac{2}{6} \approx 0.3$$

$$IQR = \frac{4}{6} \approx 0.667$$

$$b) x' = \frac{x - \mu}{\sigma} = \frac{x - 6}{2}$$

$$\text{Scaled Vector: } [-1, 1, 0, -0.5, -1.5, 0.5, 1.5]$$

$$\text{mean} = 0 \quad \text{std} = 1 \quad IQR = 2$$

$$c) x' = \frac{x - \text{median}}{\text{IQR}} = \frac{x - 6}{4}$$

Scaled Vector:

$$[-0.5, 0.5, 0, -0.25, -0.75, 0.25, 0.75]$$

median = 0, mean = 0

$$\text{std} = \frac{2}{4} = 0.5 \quad \text{IQR} = 1$$

q) Min-Max: preserves order and maps $\text{min} \rightarrow 0, \text{max} \rightarrow 1$

Standard: $\text{min} \rightarrow 0, \text{std} \rightarrow 1$

Robust: median $\rightarrow 0, \text{IQR} \rightarrow 1$, robust to outliers

3.a) Scaling before splitting leaks test information

b) The model can indirectly "see" the test dist.

c) there is optimistic bias in metrics because leakage make data easier.