

Statistical Justifications for regression

- Sample pts from unknown dist d .
- y -values are sum of unknown, non-random fn + random noise
- $\forall x_i, y_i = \underbrace{g(x_i)}_{\text{true}} + \underbrace{\varepsilon_i}_{\text{noise}} \quad \varepsilon_i \sim D', D' \text{ has } \mu=0$
- goal: find h that estimates g .
- choose $h(x) = \underbrace{\mathbb{E}_y[y|x=x]}_{\substack{\text{hope this} \\ \text{exists}}} = g(x) + \mathbb{E}[\varepsilon] = g(x)$

Least Squares \rightarrow follows max likelihood

Suppose $\varepsilon_i \sim \mathcal{N}(0, \sigma^2); y_i \sim \mathcal{N}(g(x_i), \sigma^2)$

$$\log f(y_i) = -\frac{(y_i - \mu)^2}{2\sigma^2} - \text{constant} \Leftarrow \mu = g(x_i)$$

$$\begin{aligned} l(g; X, y) &= \ln f(y_1) \dots f(y_n) = \sum_{i=1}^n \log f(y_i) \\ &= -\frac{1}{2\sigma^2} \boxed{\sum_{i=1}^n (y_i - g(x_i))^2} - \text{constant} \end{aligned}$$

max likelihood on g is equivalent to g from least squares regression

* least squares works well for normal dist. noise.

Empirical Risk

- risk is expected loss $R(h) = \mathbb{E}[L]$ over all $x \in \mathbb{R}^d$ $y \in \mathbb{R}$
- discriminative model \rightarrow don't know x 's dist
- assume x is dist with empirical dist: discrete uniform over sample pts
- Empirical risk: expected loss under empirical distribution

$$\hat{R}(h) = \frac{1}{n} \sum_{i=1}^n l(h(x_i), y_i)$$

\rightarrow why we minimize the sum of loss functions

Logistic Loss from max likelihood

what cost fn for probabilities?

prob pt. x_i in $y_i \leftarrow$ actual prob

Make β duplicate x_i 's, $y_i \beta$ in class $(1-y_i)\beta$ not

$$\mathcal{L}(h; X, y) = \prod_{i=1}^n h(x_i)^{y_i \beta} (1-h(x_i))^{(1-y_i)\beta} \leftarrow \text{almost binomial distribution}$$

$$l(h) = \beta \sum_{i=1}^n y_i \ln h(x_i) + (1-y_i) \ln(1-h(x_i))$$

$$= -\beta \sum \text{logistic loss fn}$$

\Rightarrow min all logistic losses

The bias-variance decomposition

bias: error due to inability of hypoth h to fit g
ex. fitting quadratic g with linear h

Variance: error due to fitting random noise in data
ex. fit linear g with linear h , yet $h \neq g$

$$X_i \sim D, \quad \varepsilon_i \sim D', \quad y_i = g(x_i) + \varepsilon_i \quad \text{fit } h$$

$$\text{pt. } z \in \mathbb{R}^d \quad \gamma = g(z) + \varepsilon$$

$$\mathbb{E}[\gamma] = g(z) \quad \text{Var}(\gamma) = \text{Var}(\varepsilon)$$

$$R(h) = \mathbb{E}[L(h(z), \gamma)] = \mathbb{E}[(h(z) - \gamma)^2]$$

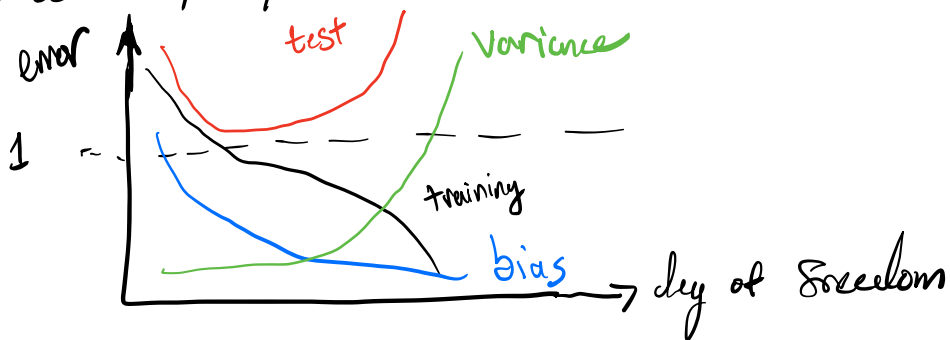
$$= \mathbb{E}[h(z)^2] + \mathbb{E}[\gamma^2] - 2 \mathbb{E}[\gamma h(z)]$$

$$= \text{Var}(h(z)) + \mathbb{E}[h(z)]^2 + \text{Var}(\gamma) + \mathbb{E}[\gamma]^2 - 2 \mathbb{E}[\gamma] \mathbb{E}[h(z)]$$

$$= \underbrace{(\mathbb{E}[h(z)] - \mathbb{E}[\gamma])^2}_{\text{bias}} + \underbrace{\text{Var}(h(z)) + \text{Var}(\gamma)}_{\text{variance}} + \underbrace{\mathbb{E}[\gamma]^2}_{\text{irreducible error}}$$

↑
"bias-variance decomposition"

- underfitting \rightarrow too much bias
- overfitting \rightarrow too much variance
- Training error reflects bias but not variance
- many dists, $\sigma \rightarrow 0$ as $n \rightarrow \infty$
- If n can be g , bias $\rightarrow 0$ " "
- \hookrightarrow it can't bias large
- good feature \downarrow bias, bad rarely increases bias
- \hookrightarrow usually increases variance
- irreducible error can't be reduced
- noise in test affects $\text{Var}(\epsilon)$
- noise in training bias/variance
- can't precisely measure bias/variance
- \hookrightarrow test w/ synthetic data



crossing point of derivative