

Question How to guarantee statistically that minimizing  $J^{(\text{train})}$  minimizes  $J^{(\text{test})}$ ?

IID Assumptions: Both training set & test set are gotten by independently randomly sampling the same probability distribution  $p_{\text{data}}$ .

Usually  $T$  is accomplished by finding some parameters  $w^*$  which minimize cost function

$J(w; E^{(\text{train})})$  —  $\left( E^{(\text{train})} \text{ means } \begin{array}{l} \text{training data} \end{array} \right)$

and measuring performance

using  $J(w^*; E^{(\text{test})})$

$\left( \text{measures the } \begin{array}{l} \text{average error} \\ \text{in } E^{(\text{train})} \end{array} \right)$

The process goes:

① Randomly sample  $p_{\text{data}}$  to get  $E^{(\text{train})}$

② Find  $w^* = \underset{w}{\operatorname{argmin}} J(w; E^{(\text{train})})$

(3) Sample  $p_{\text{data}}$  to get  $E^{(\text{test})}$

(4) Compute  $J(w^*; E^{(\text{test})})$

$$\begin{aligned}\text{So } E_{p_{\text{data}}} [J(w^*; E^{(\text{test})})] \\ &\geq E_{p_{\text{data}}} \left[ \min_w J(w; E^{(\text{test})}) \right] \\ &= E_{p_{\text{data}}} \left[ \min_w J(w; E^{(\text{train})}) \right] \text{ by IID} \\ &= E_{p_{\text{data}}} [J(w^*; E^{(\text{train})})]\end{aligned}$$

So value of test error is expected to be larger than the training error.

So performance of model depends on

(1) Minimizing training error

(2) Minimizing the gap between training error and test error

If you struggle with (1) this is a problem of "underfitting"

(2)  $\Rightarrow$  "overfitting"

(Informal) The capacity of a learning algo is the "size" of the space of functions the algo can fit. hypothesis space

Note: To make this formal, we need statistical learning theory (e.g. VC dimensions)

E.g. Linear regression with 1 variable

Case 1 no polynomial terms  
predict with  $f(x) = wx + b$  hypothesis function

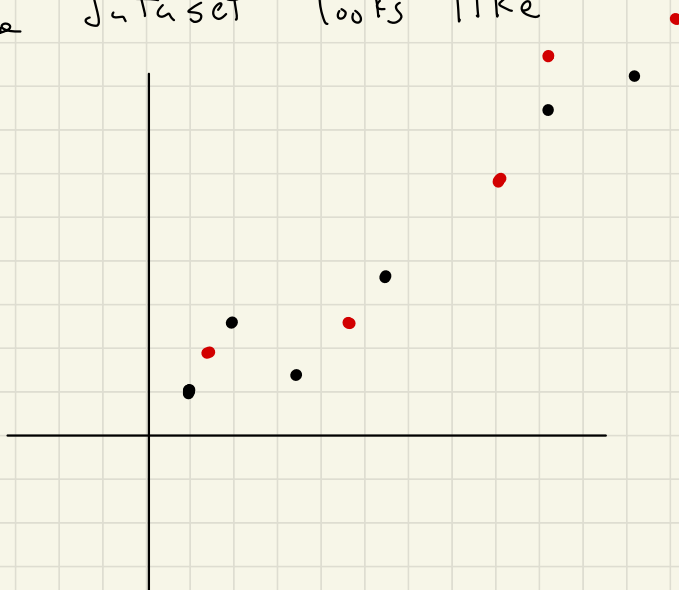
(capacity) the space of all such functions has dim 2

Case 2 poly terms up to degree 4

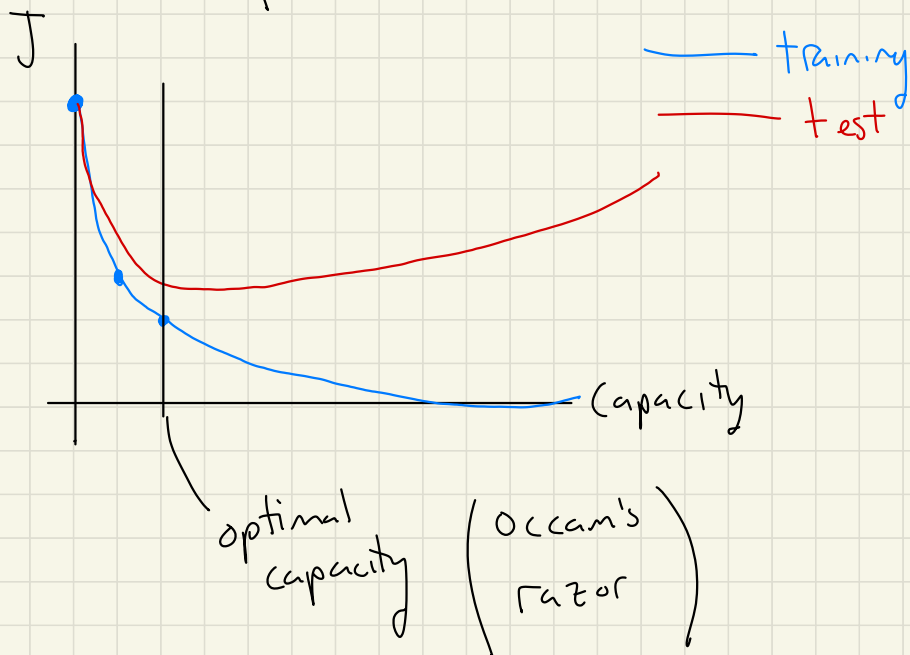
predict with  $f(x) = w_4 x^4 + w_3 x^3 + w_2 x^2 + w_1 x + b$

Now capacity is 5

Suppose dataset looks like



Let's graph (training error) versus capacity  
test error



Bayes Error (Suppose we have  $p_{\text{data}}$ )

Use  $p_{\text{data}}$  to make predictions

e.g.  $f^B(\vec{x}) = \underset{y}{\operatorname{argmax}} p_{\text{data}}(y | \vec{x})$  (supervised algo)

Such a model provides a theoretical minimal cost for any prediction model. Define

$$\text{Bayes Error} = \min_f \mathbb{E}_{p_{\text{data}}} [J(f; \vec{x}, y)]$$

where  $J(f; \vec{x}, y)$  is the cost of prediction  $f(\vec{x}) = y$

Note: Bayes Error can be nonzero

- mapping  $\vec{x}$  to  $y$  can be nondeterministic
  - $\vec{x}$  can be mapped to  $y$  in a deterministic way but using more variables than we have in  $X$ .
- 

Regularization:

with parameter  $\lambda$

Idea: Add a new term to cost function  $J$  to further restricting the hypothesis space

Example (Linear regression)

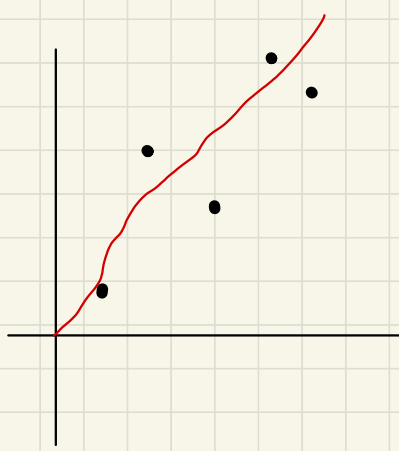
$$\text{Modifying } J(\vec{w}) = \frac{1}{n} \|\hat{X}\vec{w} - y\|_2^2 + \lambda \|\vec{w}\|_2^2$$

Since we minimize  $J$  to find  $w^*$ , this extra term promotes coefficients of  $w^*$  to be small.  $\lambda$  dictates how strongly this is promoted.

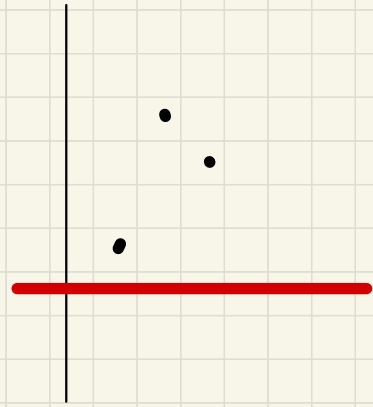
E.g. Using poly terms up to degree 10



$\lambda = 0$   
(overfit)



Somewhere in  
between  
(good fit)



$\lambda$  Large  
(underfit)

Using regularization, can start with high capacity model and "tune"  $\lambda$  to avoid overfitting/underfitting.

## Hyperparameters and Validation

Hyperparameter: Parameters which affect/control a learning algo's behavior but are not learned by the algorithm.

E.g. Polynomial regression w/ regularization

$n, \lambda$  are hyperparameters  
↑  
max degree of poly terms

Hyperparameters are usually tuned by hand to minimize test error.

Issue: When we do this, the test set is no longer a test set  
(model parameters are fit to the test set)

Solution Introduce a new "test set" called validation set

The process

Split data  $E$  into 80%  $E^{\text{train}}$ , 10%  $E^{\text{val}}$ , 10%  $E^{\text{test}}$

① Given hyperparameters  $\vec{\lambda}$  train model with  $E^{\text{train}}$  by finding  $\vec{w}_{\vec{\lambda}}^* = \underset{\vec{w}}{\operatorname{argmin}} J(\vec{w}, \vec{\lambda}; E^{\text{train}})$

② Evaluate performance by looking at  $J(\vec{w}_{\vec{\lambda}}^*; E^{\text{val}}) \leftarrow \text{(without regularization term)}$



(3) Repeat ① & ② with different choices of  $\lambda$  until find lowest value  $J(\vec{w}_{\lambda^*}^*; E^{val})$

$\longrightarrow \lambda^*$  best  $\lambda$  (tuning)

(4) Approximate generalization error by computing

$$J(\vec{w}_{\lambda^*}^*; E^{test}) \quad \left( \begin{array}{l} \text{without} \\ \text{regularization} \\ \text{term} \end{array} \right)$$

---

Tip: If  $E$  is too small, may want to  
(read about this) try k-fold cross validation  
(allows for putting more in test set)

---

in python `cross_val_score` (sklearn)

---

Bias & Variance:

Def: A point estimator for a parameter  $\theta$  associated to some distribution  $p$  is any function of random variables  $x_1, \dots, x_n$  (distributed by  $p$ )

write  $\hat{\theta} = g(x_1, \dots, x_n)$   
↑  
estimation for  $\theta$

The bias of  $\hat{\theta}$  is

$$\text{bias}(\hat{\theta}) = E_p[\hat{\theta}] - \theta$$

The variance of  $\hat{\theta}$  is

$$\text{Var}(\hat{\theta}) = E_p[(\hat{\theta} - E_p[\hat{\theta}])^2]$$

E.g.  $X_1, \dots, X_n$  IID random variables  
distributed by Bernoulli distribution  $\Omega = \{0, 1\}$

$$p(x=1) = \theta$$

$$p(x=x) = \theta^x (1-\theta)^{(1-x)}$$

Estimate  $\theta$  by  $\hat{\theta} = \frac{\#\{i | x_i = 1\}}{n} = \frac{\sum_{i=1}^n x_i}{n}$

$$\hat{\theta} = \frac{\sum_{i=1}^n x_i}{n}$$

$$\begin{aligned}
\text{bias}(\hat{\theta}) &= E_p[\hat{\theta}] - \theta \\
&= E_p\left[\frac{\sum_{i=1}^n x_i}{n}\right] - \theta \\
&= \frac{1}{n} \sum_{i=1}^n E_p[x_i] - \theta \\
&= \frac{1}{n} \sum_{i=1}^n \theta - \theta \\
&= \frac{n\theta}{n} - \theta = 0
\end{aligned}$$

Def If  $\text{bias}(\hat{\theta}) = 0$  call  $\hat{\theta}$  an unbiased estimator.

$$\begin{aligned}
\text{Var}(\hat{\theta}) &= \text{Var}\left(\frac{1}{n} \sum_{i=1}^n x_i\right) \\
&= \frac{1}{n^2} \sum_{i=1}^n V(x_i) \\
&= \frac{1}{n^2} \sum_{i=1}^n \theta(1-\theta) \\
&= \frac{\theta(1-\theta)}{n} \longrightarrow 0 \quad \text{as } n \rightarrow \infty
\end{aligned}$$

Exercises ① Gaussian distribution  $\mathcal{N}(x; \mu, \sigma^2)$

estimate  $\mu$  by IID random variables  $X_1, \dots, X_n$

$$\hat{\mu} = \frac{X_1 + \dots + X_n}{n}$$

Show  $\text{Bias}(\hat{\mu}) = 0$

② Estimate  $\sigma^2$  by  $\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \hat{\mu})^2$

$$\text{Show } \text{Bias}(\hat{\sigma}^2) = -\frac{\sigma^2}{n}$$

$$\text{Var}(\hat{\sigma}^2) = ?$$

③ Estimate  $\sigma^2$  by  $\hat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \hat{\mu})^2$

$$\text{Show } \text{bias}(\hat{\sigma}^2) = 0 \quad \text{Var}(\hat{\sigma}^2) = ?$$

(called the unbiased variance estimator)

Def: The mean squared error (MSE) of  $\hat{\theta}$

$$\text{is } \text{MSE}(\hat{\theta}) = E_{\mathbf{p}}[(\hat{\theta} - \theta)^2]$$

Lemma  $MSE(\hat{\theta}) = Var(\hat{\theta}) + Bias(\hat{\theta})^2$

Proof:  $MSE(\hat{\theta}) = E_p[(\hat{\theta} - \theta)^2]$

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

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

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

Note  $E[\hat{\theta}] - \theta$  is constant

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

$$+ 2(E[\hat{\theta}] - \theta) \underbrace{E[(\hat{\theta} - E[\hat{\theta}])]}_{E[\hat{\theta}] - E[\hat{\theta}] = 0}$$

$$E[\hat{\theta}] - E[\hat{\theta}] = 0$$

$$= Var(\hat{\theta}) + bias(\hat{\theta})^2$$



# Maximum Likelihood Estimators

Back to usual setup:  $p_{\text{data}}$  generating

$$X = X^{\text{train}}, X^{\text{test}} \quad \left[ \begin{array}{l} X \text{ has rows } \vec{x}^{(1)}, \dots, \vec{x}^{(m)} \\ \text{treated as random} \\ \text{variables distributed} \\ \text{by } p_{\text{data}} \end{array} \right]$$

Goal: Approximate  $p_{\text{data}}(\vec{x})$ . Let

$p_{\text{model}}(\vec{x}, \vec{\theta})$  be an approximation for  $p_{\text{data}}$  depending on parameter  $\vec{\theta}$

E.g. If we expect  $p_{\text{data}}$  to be gaussian

$$p_{\text{model}}(\vec{x}, \theta) = \mathcal{N}(\vec{x}; \vec{\mu}_0, \Sigma_0)$$

Def The maximum likelihood estimator for  $\vec{\theta}$  is

$$\bar{\theta}_{ML} = \underset{\bar{\theta}}{\operatorname{argmax}} p_{\text{model}}(\bar{x}^{(1)}, \dots, \bar{x}^{(m)}; \bar{\theta})$$