

$$P(x|y), \quad P(y)$$

Maximum Likelihood Estimation

Find parameters $P(x|y)$ and $P(y)$ that maximize the probability of observing the training set.

MAP rule: Once parameters are discovered, on a new example (x)

$$p(y=0|x) = \frac{p(x|y=0) \cdot p(y=0)}{p(x|y=0) \cdot p(y=0) + p(x|y=1) \cdot p(y=1)}$$

$$p(y=1|x) = \frac{p(x|y=1) \cdot p(y=1)}{p(x)}$$

Output 1 if $p(y=1|x) > p(y=0|x)$
0 otherwise

Gaussian Discriminant Analysis

$$p(x|y) \sim N(\mu_y, \Sigma) \leftarrow$$

$$p(y) \sim \text{Ber}(\phi)$$

$$\text{MAP rule: } \log \left[\frac{p(y=1|x)}{p(y=0|x)} \right] > 0$$

$$P(y=0|x)$$

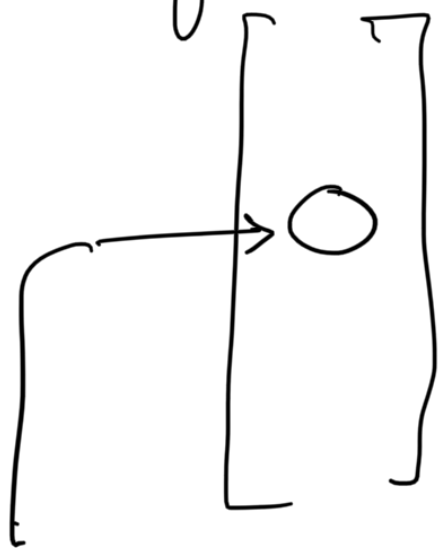
⇒ linear classifier.

$$P(y=1|x) = \frac{1}{1 + e^{-\theta^T x}}$$

Naive Bayes:

Classification: Given an email
classify if it is spam.

Input features: Given email



← One position for
every word in
same vocabulary

0 if this word does not appear in the
email

1 if this word appears in email.

Discriminative Learning Algorithm

$P(y=1|x) \leftarrow$ learn

\hookrightarrow Bern (ϕ_x)

parameters is 2^d ($x \in \mathbb{R}^d$)

Generative Learning Algorithm

$$\text{parameters: } p(y=1) = \phi_y$$

$$p(x|y=0)$$

$$p(x|y=1)$$

Naive Bayes assumption:

Probability that the j th word appears in an email (spam/non-spam) is independent of the k th word appearing in an email.

$$p(x|y=0) = \prod_{j=1}^d p(x_j=1|y=0)$$

$$p(x_j=1|y=0) = \phi_{j|0} \quad \leftarrow 2d$$

$$p(x_j=1|y=1) = \phi_{j|1} \quad \leftarrow 2d$$

$$\text{parameters: } p(y=i) = \phi_y$$

$$\text{if } y=i \in \{0,1\} \quad p(x_j=1|y=i) = \phi_{j|i}$$

$$L(\phi_y, \phi_{j|0}, \phi_{j|1})$$

$$= \prod_{i=1}^n \prod_{j=1}^d \phi_{j|y}^{1[x_j^{(i)}=1]} (1-\phi_{j|y})^{1[x_j^{(i)}=0]}$$

$\prod_{y=1}^n \phi_y^{1[y=1]} (1-\phi_y)^{1[y=0]}$

$\Psi_y \quad (y)$

We will maximize the log of the likelihood

$$\phi_y = \sum_{i=1}^n 1[y^{(i)} = 1]$$

$$\phi_{j|y=0} = \frac{\sum_{i=1}^n 1[x_j^{(i)} = 1 \wedge y^{(i)} = 0]}{\sum_{i=1}^n 1[y^{(i)} = 0]} \leftarrow$$

$$\phi_{j|y=1} = \frac{\sum_{i=1}^n 1[x_j^{(i)} = 1 \wedge y^{(i)} = 1]}{\sum_{i=1}^n 1[y^{(i)} = 1]}$$

$$\phi_{j=\text{COVID}|0} = 0 \leftarrow$$

$$\phi_{j=\text{COVID}|1} = 0$$

$$\begin{aligned} \rightarrow P(y=1|x) &= \frac{p(x|y=1) \cdot p(y=1)}{p(x|y=0) \cdot p(y=0) + p(x|y=1) \cdot p(y=1)} \\ &= \frac{0}{\frac{0}{1} \phi_{j||} \phi_y + 0} \end{aligned}$$

$$= \frac{0}{0}$$

1 || something.

Laplace Smoothing

Random Variable $Z \sim \text{Multinomial}(k)$

$$Z \in \{1, 2, \dots, k\}, \phi_i = P[Z=i]$$

$$\phi_j = \frac{1 + \sum_{i=1}^k 1[Z=i]}{k+n}$$

Back to email classification

$$\phi_j | y=s = \frac{1 + \sum 1[x_j^{(i)}=1 \wedge y^{(i)}=s]}{2 + \sum 1[y^{(i)}=s]}$$

Generalization:

Goal of learning: Discover a model of the world that explains data.

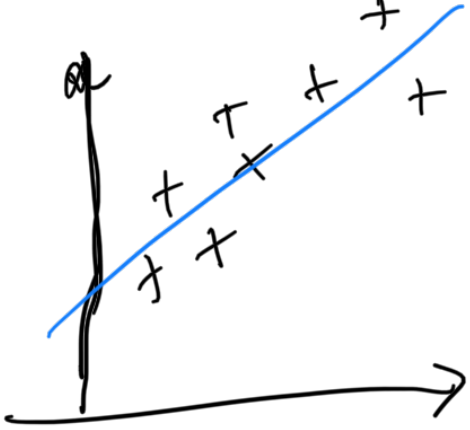
Minimizing training is a way to try to ensure that our learnt model will do well on new test data.

Assuming: There is distribution D and training set is generated by sampling from D and new examples

will come from D .

Real goal is minimize

$$E_{x \sim D} [\text{loss}]$$



Linear

High Training error

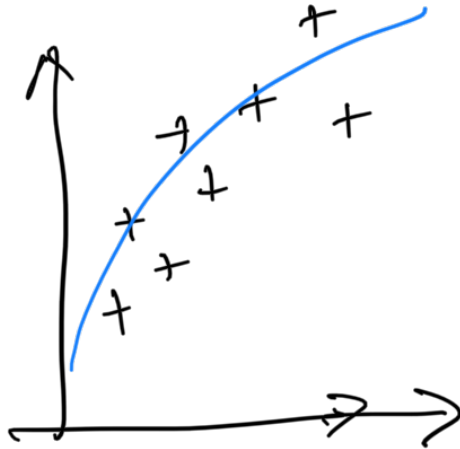
Regression

High Bias

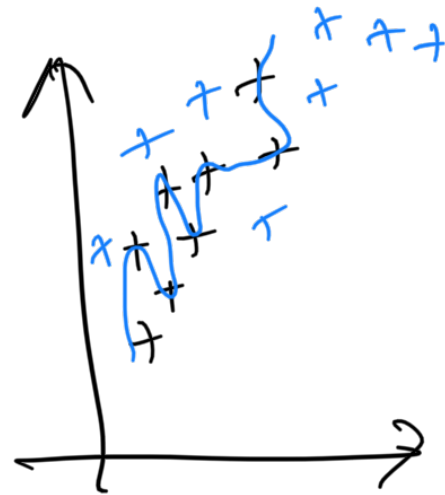
Underfitting - [High training + test error]

Cannot be addressed by more training data

Bias / Variance



Quadratic



High poly

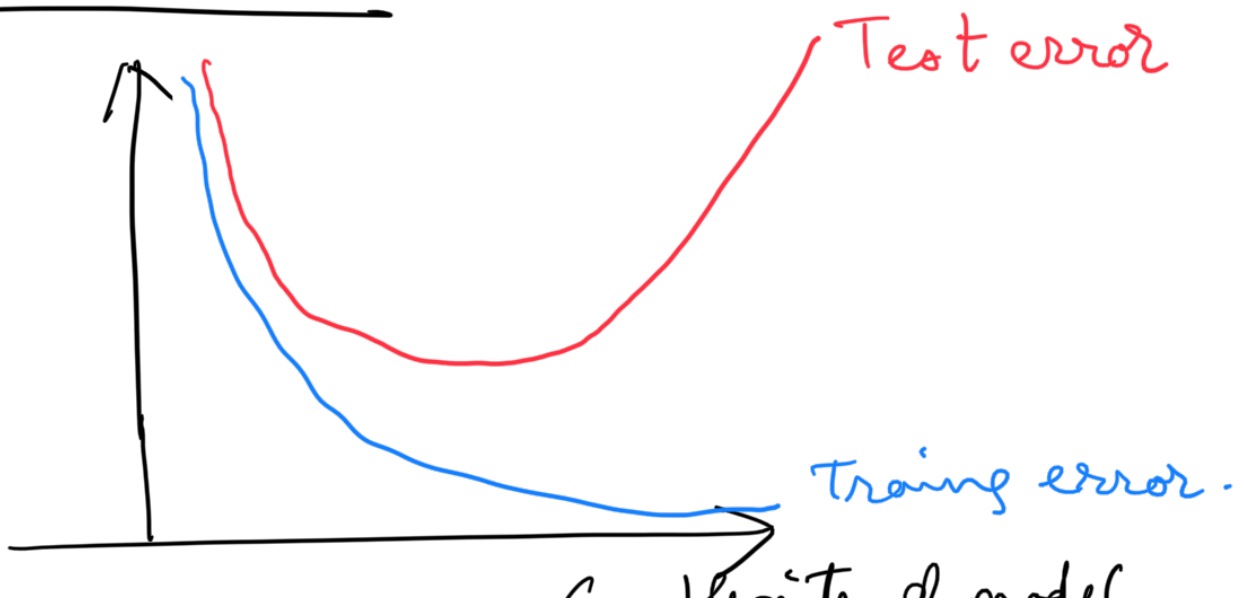
Lowest training error

High Variance

Overfitting.

Low training error
High test error

Addressed increasing training data.



Complexity of model