

Unbalanced Datasets

- For calculating accuracy we have implicitly assumed that there are the same number of positive and negative examples in the dataset (which is known as a balanced dataset).

balanced accuracy (BA)

$$BA = \frac{TPR + TNR}{2}$$

- It doesn't hold good always

Matthews correlation coefficient (MCC)

$$MCC = \frac{TP \times TN - FP \times FN}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}}$$

Example - MCC

- Given a sample of 13 pictures, 8 of cats and 5 of dogs, where cats belong to class 1 and dogs belong to class 0,
- Actual = [1,1,1,1,1,1,1,1,0,0,0,0,0],
- Prediction = [0,0,0,1,1,1,1,1,0,0,0,1,1]
- Classifier makes 8 accurate predictions and misses 5: 3 cats wrongly predicted as dogs (first 3 predictions) and 2 dogs wrongly predicted as cats (last 2 predictions).

		Actual class	
		Cat	Dog
Predicted class	Cat	5	2
	Dog	3	3

		Actual class	
		P	N
Predicted class	P	TP	FP
	N	FN	TN

Example - MCC

- $MCC = [(5*3) - (2*3)] / \sqrt{[(5+2)*(5+3)*(3+2)*(3+3)]} = 9 / \sqrt{1680} = 0.219$

- Bayes Rule

$$P(C_i|X_j) = \frac{P(X_j|C_i)P(C_i)}{P(X_j)}$$

Example – Multiclass Bayesian classification

- What to do in the evening based on whether you have an assignment deadline and what is happening.

Deadline?	Is there a party?	Lazy?	Activity
Urgent	Yes	Yes	Party
Urgent	No	Yes	Study
Near	Yes	Yes	Party
None	Yes	No	Party
None	No	Yes	Pub
None	Yes	No	Party
Near	No	No	Study
Near	No	Yes	TV
Near	Yes	Yes	Party
Urgent	No	No	Study

New Test Data: Suppose that you have deadlines looming, but none of them are particularly urgent, that there is no party on, and that you are currently lazy.

- $P(\text{Party}) \times P(\text{Near} \mid \text{Party}) \times P(\text{No Party} \mid \text{Party}) \times P(\text{Lazy} \mid \text{Party})$
- $P(\text{Study}) \times P(\text{Near} \mid \text{Study}) \times P(\text{No Party} \mid \text{Study}) \times P(\text{Lazy} \mid \text{Study})$
- $P(\text{Pub}) \times P(\text{Near} \mid \text{Pub}) \times P(\text{No Party} \mid \text{Pub}) \times P(\text{Lazy} \mid \text{Pub})$
- $P(\text{TV}) \times P(\text{Near} \mid \text{TV}) \times P(\text{No Party} \mid \text{TV}) \times P(\text{Lazy} \mid \text{TV})$

Using the data above these evaluate to:

$$\begin{aligned} P(\text{Party} \mid \text{near (not urgent) deadline, no party, lazy}) &= \frac{5}{10} \times \frac{2}{5} \times \frac{0}{5} \times \frac{3}{5} \\ &= 0 \end{aligned}$$

$$\begin{aligned} P(\text{Study} \mid \text{near (not urgent) deadline, no party, lazy}) &= \frac{3}{10} \times \frac{1}{3} \times \frac{3}{3} \times \frac{1}{3} \\ &= \frac{1}{30} \end{aligned}$$

$$\begin{aligned} P(\text{Pub} \mid \text{near (not urgent) deadline, no party, lazy}) &= \frac{1}{10} \times \frac{0}{1} \times \frac{1}{1} \times \frac{1}{1} \\ &= 0 \end{aligned}$$

$$\begin{aligned} P(\text{TV} \mid \text{near (not urgent) deadline, no party, lazy}) &= \frac{1}{10} \times \frac{1}{1} \times \frac{1}{1} \times \frac{1}{1} \\ &= \frac{1}{10} \end{aligned}$$

So based on this you will be watching TV tonight.

Basic statistics

- Averages – mean, median, mode
- Variance
- The variance of the set of numbers is a measure of how spread out the values are. It is computed as the sum of the squared distances between each element in the set and the expected value of the set (the mean, μ)

$$\text{var}(\{\mathbf{x}_i\}) = \sigma^2(\{\mathbf{x}_i\}) = E((\{\mathbf{x}_i\} - \mu)^2) = \sum_{i=1}^N (\mathbf{x}_i - \mu)^2$$

- Standard deviation
- <https://www.itl.nist.gov/div898/handbook/pmc/section5/pmc541.htm>

Basic statistics

- Covariance
- If two variables are independent, then the covariance is 0 (the variables are then known as uncorrelated), while if they both increase and decrease at the same time, then the covariance is positive, and if one goes up while the other goes down, then the covariance is negative.

$$\text{cov}(\{x_i\}, \{y_i\}) = E(\{x_i\} - \mu)E(\{y_i\} - \nu)$$

Basic statistics

Covariance matrix

The covariance can be used to look at the correlation between all pairs of variables within a set of data. We need to compute the covariance of each pair, and these are then put together into what is imaginatively known as the covariance matrix.

$$\Sigma = \begin{pmatrix} E[(x_1 - \mu_1)(x_1 - \mu_1)] & E[(x_1 - \mu_1)(x_2 - \mu_2)] & \dots & E[(x_1 - \mu_1)(x_n - \mu_n)] \\ E[(x_2 - \mu_2)(x_1 - \mu_1)] & E[(x_2 - \mu_2)(x_2 - \mu_2)] & \dots & E[(x_2 - \mu_2)(x_n - \mu_n)] \\ \dots & \dots & \dots & \dots \\ E[(x_n - \mu_n)(x_1 - \mu_1)] & E[(x_n - \mu_n)(x_2 - \mu_2)] & \dots & E[(x_n - \mu_n)(x_n - \mu_n)] \end{pmatrix}$$

Example

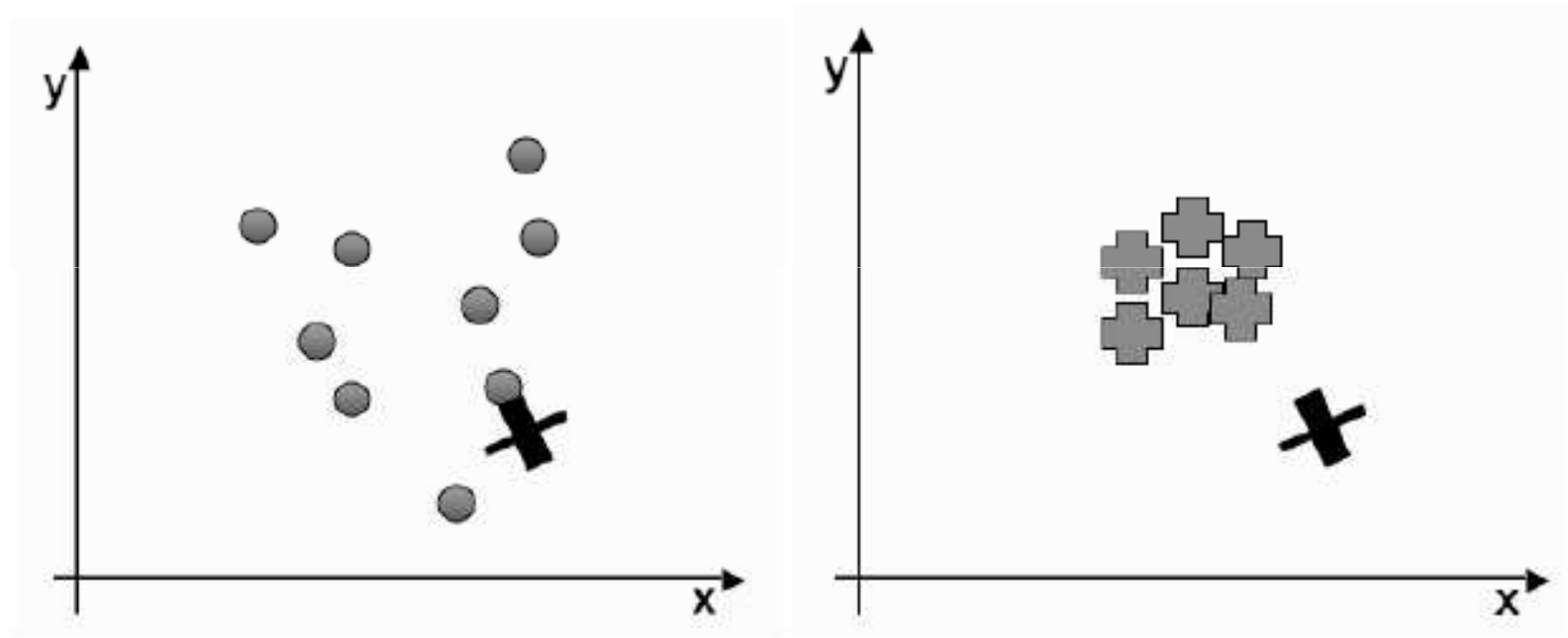


FIGURE 2.13 Two different datasets and a test point.

Gaussian Distribution

- One dimensional / multi dimensional

$$p(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

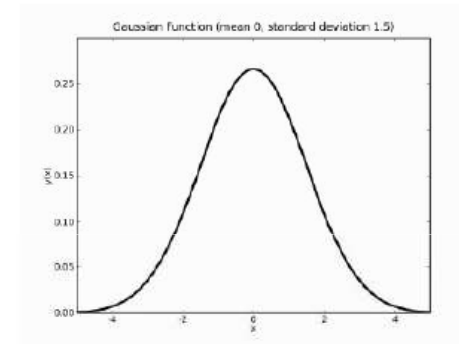


FIGURE 2.14 Plot of the one-dimensional Gaussian curve.

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu)\right)$$

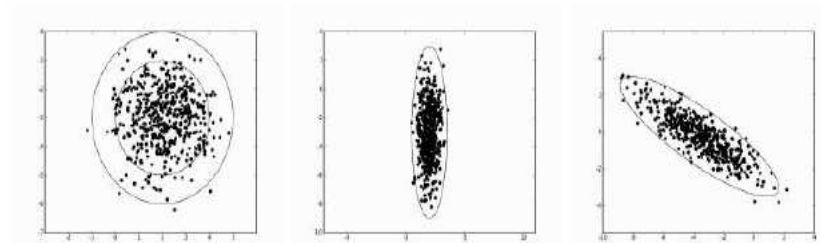


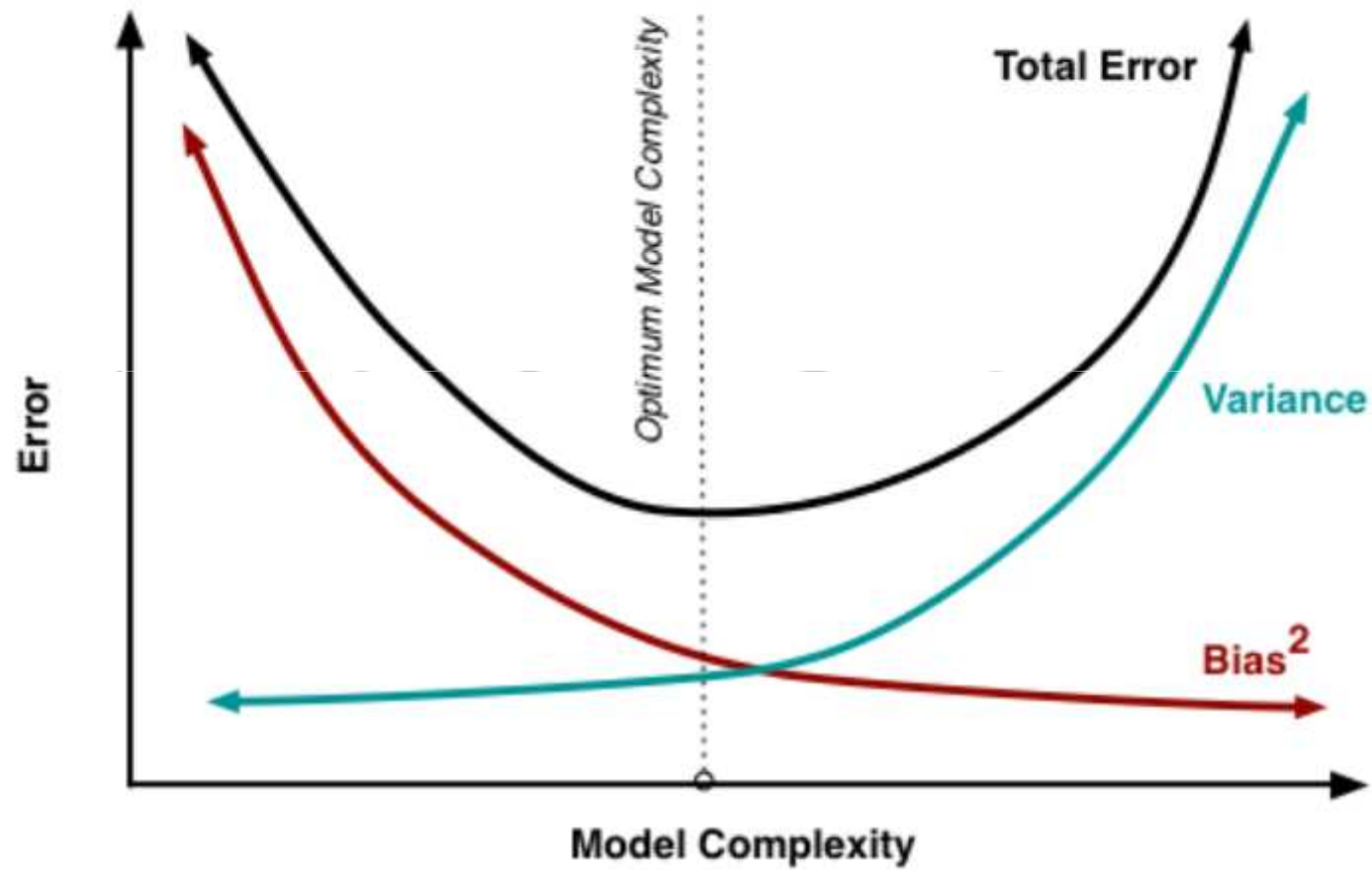
FIGURE 2.15 The two-dimensional Gaussian when (*left*) the covariance matrix is the identity, (*centre*) the covariance matrix has elements on the leading diagonal only, and (*right*) the general case.

Bias Variance Tradeoff



- Train ML – make choice about the model, fit parameters
- More degrees of freedom the algorithm has, the more complicated the model that can be fitted - overfitting, and requiring more training data, need for validation data
- There is another way to understand this idea that more complex models do not necessarily result in better results. Some people call it the **bias-variance dilemma** rather than a **tradeoff**.
- A model can be bad for **two different reasons**.
 - Either not accurate and doesn't match the data well -> bias
 - or it is not very precise and there is a lot of variation in the results -> statistical variance.

- Eg: Consider the difference between a straight line fit to some data and a high degree polynomial, which can go precisely through the datapoints.
- The straight line has no variance at all, but high bias since it is a bad fit to the data in general.
- The spline can fit the training data to arbitrary accuracy, but the variance will increase.
- Note that the variance probably increases by rather less than the bias decreases, since we expect that the
- spline will give a better fit.
- Some models are definitely better than others, but choosing **the complexity of the model is important** for getting good results.
- The variance - σ^2 changes depending on the particular training set that was used, while the bias - average error
- Expectation of sum of squares error = $\text{noise}^2 + \text{variance} + \text{bias}^2$



Bias and variance contributing to total error.

