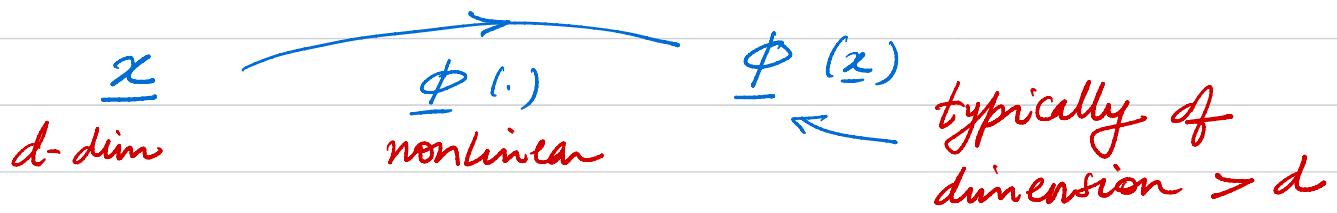
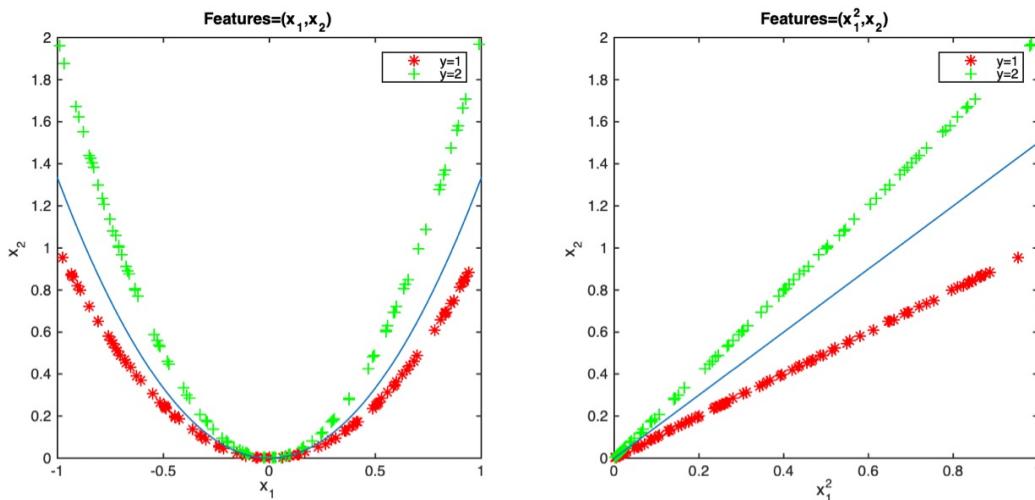


Kernel Trick

Map feature vector to new space



why? Easier to separate training data using linear boundaries in new space

Dot Products and Kernels

- Building linear classifiers often involves only dot (inner) products between vectors
 - Classifier boundary : $\underline{w}^T \underline{x} = \beta$
 - Euclidean distance : $\Delta(\underline{x}_1, \underline{x}_2) = \sqrt{(\underline{x}_1 - \underline{x}_2)^T (\underline{x}_1 - \underline{x}_2)}$
 - If we map \underline{x} to $\phi(\underline{x})$, then dot products in new space formed using
- $$\underline{\phi(\underline{x})}^T \underline{\phi(\underline{r})} \triangleq k(\underline{x}, \underline{r}) \leftarrow \text{kernel}$$

Mercer Kernel If κ satisfies the following :

- (1) $\kappa(\cdot, \cdot)$ is continuous,
- (2) $\kappa(\cdot, \cdot)$ is symmetric, i.e., $\kappa(\underline{x}, \underline{r}) = \kappa(\underline{r}, \underline{x})$
- (3) $\kappa(\cdot, \cdot)$ is positive semi-definite.

Then κ is called a Mercer kernel, and for any Mercer kernel it can be shown that there exists $\underline{\phi}$ s.t. $\kappa(\underline{x}, \underline{r}) = \underline{\phi}(\underline{x})^T \underline{\phi}(\underline{r})$

Examples of kernels

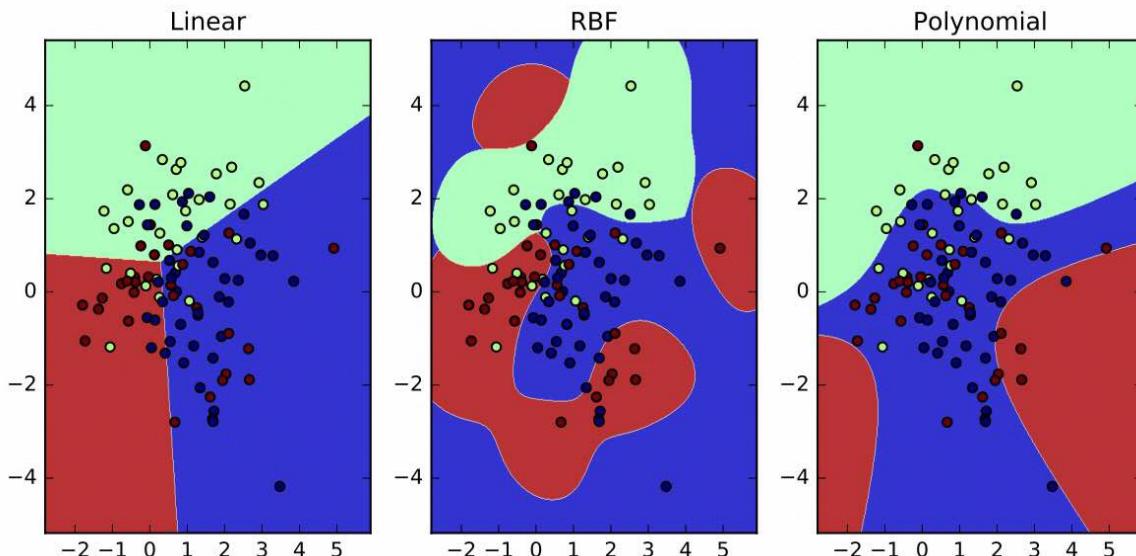
1. Linear : $\kappa(\underline{x}, \underline{r}) = c \underline{x}^T \underline{r}$

2. Radial Basis Function (RBF) : $\kappa(\underline{x}, \underline{r}) = h(\|\underline{x} - \underline{r}\|)$

e.g. Gaussian RBF : $\kappa(\underline{x}, \underline{r}) = e^{-c \|\underline{x} - \underline{r}\|^2}$ ↑ monotonic

3. Polynomial kernel : $\kappa(\underline{x}, \underline{r}) = (1 + \underline{x}^T \underline{r})^l$, $l=1, 2, \dots$

4. Sigmoid : $\kappa(\underline{x}, \underline{r}) = \tanh(c \underline{x}^T \underline{r} + b)$ not Mercer



SVM

How to Handle Labelled Data

- Need to be careful not to use labelled data in way that leads to **circular conclusions**
- Classifier that works perfectly on training data, i.e., $\text{Err}_{\text{train}} = 0$, may perform poorly on actual inputs - **poor generalization**

$$\text{Recall: } f(\tilde{x}) = \begin{cases} \tilde{y} & \text{if } (\tilde{x}, \tilde{y}) \in T \\ -1 & \text{else} \end{cases}$$

- Two fundamental steps in designing good $f(\cdot)$:

1. Model Selection Given several models, estimate their performance to choose "best" one
2. Model Assessment Having decided on best model, estimate how well it generalizes. $\hat{\text{Err}}_{\text{pred}}$

Split labelled data into 3 parts:



CAUTION Don't touch test data until **Test stage**

Model Selection:

- Use training data on several models ; e.g. k-NN, LDA,..
- Use Validation data to compute Err_{val} for models
- Select best - smallest Err_{val}

Cross-Validation

- In many applications, limited labelled data
- Cannot afford to split into 3 parts and use exclusively for training, validation, and testing

Step 0 Set aside test data

Step 1 split remaining labelled data into K approximately equal parts $\tau_1, \tau_2, \dots, \tau_K$

τ_1	τ_2	\dots	τ_K
----------	----------	---------	----------

τ_e - called the e -th fold

Step 2 For $e=1, 2, \dots, K$,

- Use all data except that in τ_e to train models
- Use τ_e to validate models

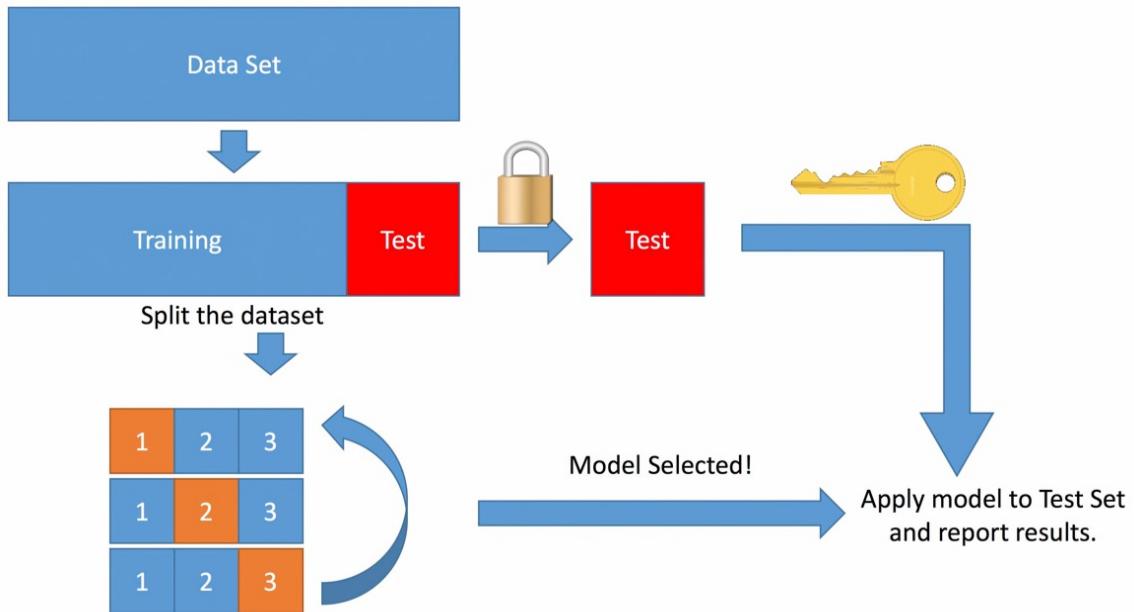
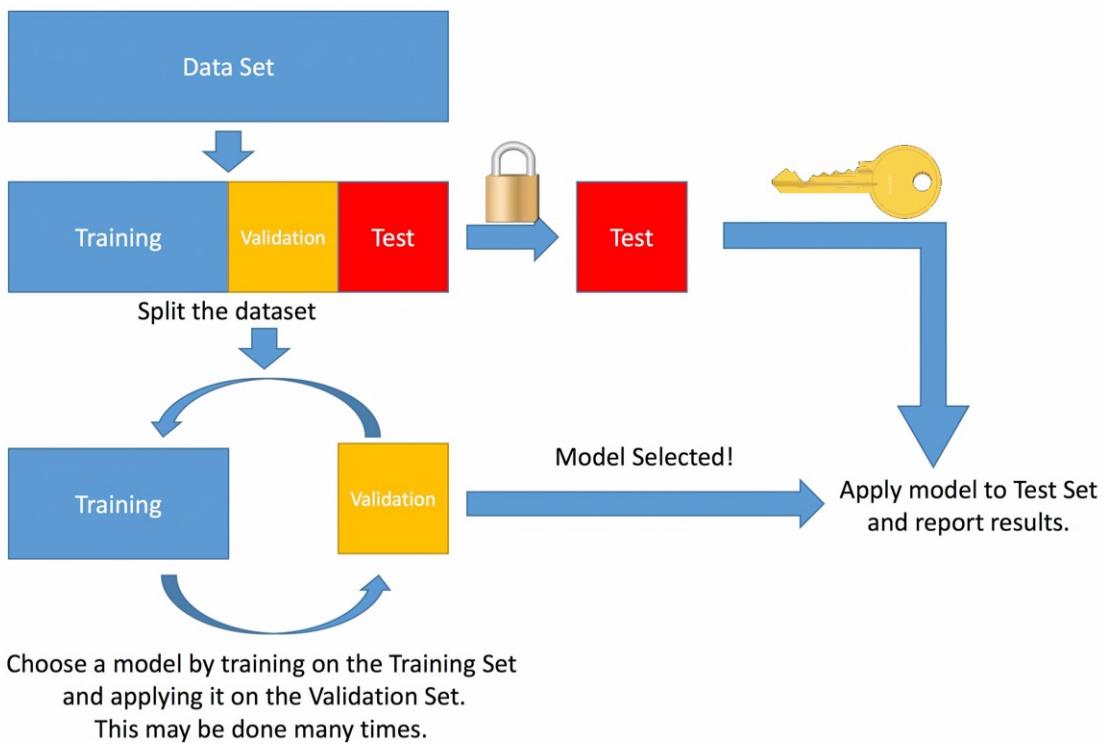
L_e = validation error for a given model

Step 3 Compute cross-validation error:

$$\text{Err}_{cv} = \frac{1}{K} \sum_{e=1}^K L_e .$$

Step 4 Use Err_{cv} of different models to select best

Step 5 Evaluate error of selected model on test data. This gives an accurate estimate of Err_{pred} for selected model



Choose a model applying cross-validation using the training set. This may be done many times.